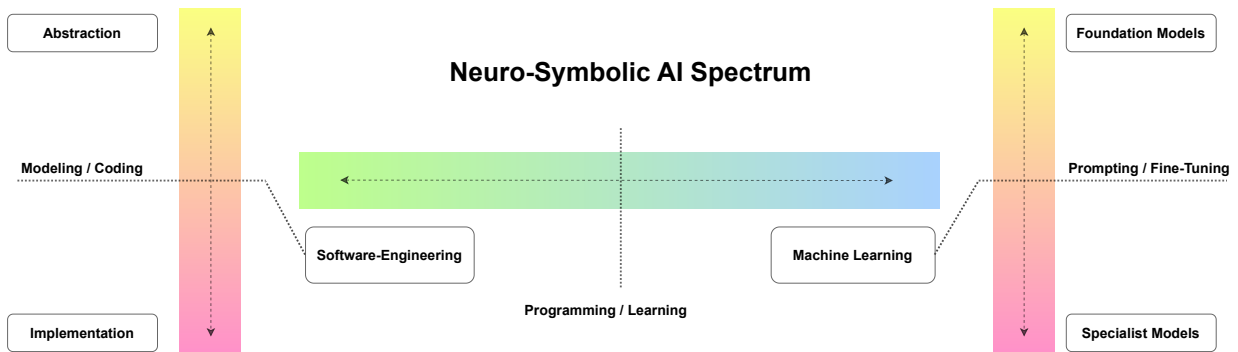


# SYMBOLICAI: A FRAMEWORK FOR LOGIC-BASED APPROACHES COMBINING GENERATIVE MODELS AND SOLVERS

Marius-Constantin Dinu<sup>\* † ‡</sup>    Claudiu Leoveanu-Condrei<sup>† ‖</sup>    Markus Holzleitner<sup>‡</sup>  
 Werner Zellinger<sup>‡ §</sup>    Sepp Hochreiter<sup>‡</sup>  
 ExtensityAI<sup>†</sup>    Johannes Kepler University<sup>‡</sup>    RICAM<sup>§</sup>    Amazon Devices<sup>‖</sup>

## ABSTRACT

We introduce *SymbolicAI*, a versatile and modular framework employing a logic-based approach to concept learning and flow management in generative processes. SymbolicAI enables the seamless integration of generative models with a diverse range of solvers by treating large language models (LLMs) as semantic parsers that execute tasks based on both natural and formal language instructions, thus bridging the gap between symbolic reasoning and generative AI. We leverage probabilistic programming principles to tackle complex tasks, and utilize differentiable and classical programming paradigms with their respective strengths. The framework introduces a set of polymorphic, compositional, and self-referential operations for multi-modal data that connects multi-step generative processes and aligns their outputs with user objectives in complex workflows. As a result, we can transition between the capabilities of various foundation models with in-context learning capabilities and specialized, fine-tuned models or solvers proficient in addressing specific problems. Through these operations based on in-context learning our framework enables the creation and evaluation of explainable computational graphs. Finally, we introduce a quality measure and its empirical score for evaluating these computational graphs, and propose a benchmark that compares various state-of-the-art LLMs across a set of complex workflows. We refer to the empirical score as the "Vector Embedding for Relational Trajectory Evaluation through Cross-similarity", or *VERTEX* score for short. The [framework codebase](#)<sup>1</sup> and [benchmark](#)<sup>2</sup> are linked below.



**Figure 1:** Our neuro-symbolic framework enables a seamless transition between symbolic and differentiable programming, each with distinct dynamics and strengths. Differentiable programming provides access to foundational and specialist models. Classical programming, on the other hand, shifts between abstraction and implementation, focusing on high-level concepts before delving into the details of implementation.

<sup>1</sup> SymbolicAI framework: <https://github.com/ExtensityAI/symbolicai>

<sup>2</sup> Evaluation benchmark: <https://github.com/ExtensityAI/benchmark>

\* Correspondence to: [dinu@ml.jku.at](mailto:dinu@ml.jku.at), {marinus, leo}@extensity.ai

‖ Work done outside of Amazon.

## 1 INTRODUCTION

The recent surge in generative AI, particularly involving large language models (LLMs), has demonstrated their wide-ranging applicability across various domains (Badita, 2022; Degraeve, 2022). These models have enhanced the functionality of tools for search-based interactions (YouWrite, 2022; Writesonic, 2022; Microsoft, 2023), program synthesis (Jain et al., 2021; Romera-Paredes et al., 2023; Key et al., 2023), chat-based interactions (ReplikaAI, 2016; OpenAI, 2022; Google, 2023), and many more. Moreover, language-based approaches have facilitated connections between different modalities, enabling text-to-image (Ramesh et al., 2021; Saharia et al., 2022), text-to-video (Singer et al., 2022), text-to-3D (Poole et al., 2022), text-to-audio (Oord et al., 2016; Wang et al., 2017), and text-to-code (Wang et al., 2021b; Lu et al., 2021; Li et al., 2022b) transformations, to name a few. Consequently, by training on vast quantities of unlabelled textual data, LLMs have been shown to not only store factual knowledge (Petroni et al., 2019; Kassner et al., 2020) and approximate users’ intentions to some extent (Andreas, 2022), but also to unlock deep specialist capabilities through innovative prompting techniques (Nori et al., 2023).

Despite their versatility, current LLMs face challenges such as fallacious reasoning and the generation of erroneous content, commonly referred to as hallucinations (Jones & Steinhardt, 2022). These limitations highlight the importance of integrating complementary symbolic methods to validate and guide the generative processes of LLMs, ensuring more accurate and reliable outputs. In parallel, efforts have focused on developing tool-based approaches (Schick et al., 2023) or template frameworks (Chase, 2023) to extend LLMs’ capabilities and enable a broader spectrum of applications. However, these efforts only partially capture the potential inherent in leveraging LLMs as *semantic parsers*. In contrast to parsers for structured languages a semantic parser is able to break down unstructured human language into semantically meaningful components and transform those into a structured form. While traditionally semantic parsing has been a role filled by specialized algorithms and models, we posit that LLMs, through their training on diverse linguistic data, have developed the ability to perform semantic parsing as part of their broader natural language processing capabilities. In turn, we identify LLMs as a central component in creating sophisticated neuro-symbolic (NeSy) AI systems. These systems integrate symbolic and sub-symbolic concepts and utilize the capabilities of semantic parsing to develop symbolic expressions that enable new probabilistic programming paradigms.

We introduce *SymbolicAI*, a compositional NeSy framework able to represent and manipulate multi-modal and self-referential structures (Schmidhuber, 2007; Fernando et al., 2023). SymbolicAI augments the generative process of LLMs with in-context learning operations, realized through functional primitives, and enables the creation of versatile applications through in-context learning (Wei et al., 2022a). These operations enable logic-based components that guide the generative process and enable a modular NeSy system, including a wide range of existing solvers, formal language engines for mathematical expression evaluation, theorem provers, knowledge bases, and search engines for information retrieval. SymbolicAI exposes these solvers as building blocks for constructing compositional functions as computational graphs, making it possible to bridge classical and differentiable programming paradigms with the aim to create *domain-invariant problem solvers*. In designing the architecture of SymbolicAI, we drew inspiration from a body of evidence that suggests the human brain possesses a selective language processing module (Macswaney, 2002; Fedorenko et al., 2010; Menenti et al., 2011; Regev et al., 2013; Scott et al., 2016; Deniz et al., 2019; Hu et al., 2022), prior research on cognitive architectures (Newell & Simon, 1956; Newell et al., 1957; Newell & Simon, 1972; Newell, 1990; Laird, 2022), and the significance of language on the structure of semantic maps in the human brain (Huth et al., 2016). We consider language as a central processing module, distinct from other cognitive processes such as reasoning or memory (Paischer et al., 2022; 2023). We hypothesize that such a central processing module based in language is a core component of broad AI systems (see Appendix Section A) and enables the development of fully autonomous AI systems for decision-making.

A significant challenge encountered in the development of our framework pertained to the evaluation of LLMs when used as semantic parsers in a NeSy workflow. Current evaluation of generated content relies on metrics for single-step generative processes, such as the BLEU score (Papineni et al., 2002). These metrics are not suitable for evaluating multi-step generative processes. BLEU has limitations, as it measures  $n$ -gram-based overlap of generated output with a reference that does not consider the semantic meaning. As a result, BLEU fails to capture semantic equivalence, especially in more complex tasks. More recent metrics such as CIDEr (Vedantam et al., 2014) or SPICE (Anderson et al., 2016) are also not suitable for our problem, either because they are built upon BLEU (in case of CIDEr) or designed with inductive biases specifically for image captioning.

Therefore, alongside our framework we introduce a quality measure (and its empirical score) for multi-step generative processes based on semantic meaning. We refer to our score as "Vector Embedding for Relational Trajectory Evaluation through Cross-similarity", or VERTEX score for short. Our VERTEX score uses embeddings to compare node distributions within a computational graph. It measures the semantic meaning across the distributional path by computing at each node the cross-similarity between the generated embeddings and embeddings sampled from a reference

distribution. Furthermore, the VERTEX score is designed such that it can be used as a reward signal in a reinforcement learning setting (Sutton, 1984). Finally, we propose a benchmark for evaluating complex workflows. We define a set of basic evaluations, particularly associative predictions based on in-context learning, multi-modal bindings for tool utilization, and program synthesis for subroutine execution. Furthermore, we introduce complex evaluations for logic-based components and hierarchical computational graphs.

In summary, the key contributions presented in this work are as follows:

- We introduce SymbolicAI, a logic-based framework for concept learning and flow management in generative processes, enabling seamless integration with a wide range of foundation models and solvers.
- We leverage LLMs as semantic parsers to enable the creation of complex computational graphs by combining symbolic expressions with probabilistic programming paradigms.
- We introduce a quality measure and its empirical score alongside a benchmark designed for multi-step generative processes for comparing LLMs across a wide range of complex tasks.

## 2 RELATED WORK

**Symbolic Methods** The field of symbolic AI has its foundations in the works of the Logic Theorist (LT) (Newell & Simon, 1956) and the General Problem Solver (GPS) (Newell et al., 1957). These programs represented the first steps towards automated reasoning and problem-solving utilizing symbolic representations. Despite their advancements, both faced challenges in dealing with the complexity of real-world problems, particularly due to the combinatorial nature of the solution space. To address these limitations, the Soar (Laird et al., 1987) cognitive architecture was developed, advancing the notion that intelligent behavior results from goal-oriented search through a problem space (Newell & Simon, 1972; McCarthy et al., 2006), with each step consisting of selecting and applying operators. Soar introduced components like reinforcement learning, impasses, sub-states, and chunking to enhance its problem-solving capabilities. It also demonstrated the importance of learning from experiences to adapt and improve performance over time. However, Santoro et al. (2022) emphasizes the subjectivity of symbols and suggests that human-like symbolic fluency could develop in machines through learning algorithms immersed in socio-cultural contexts. This perspective, anchored in the notion that symbols are triadic and their meaning emerges from consensus, seeks to move away from traditional symbolic AI methodologies towards AI that adaptively learns meaning and behaviors from human-like experiences. The goal is to cultivate machines that demonstrate symbolic behaviors across a spectrum of competencies, potentially mirroring the evolutionary and social learning processes observed in humans. Lastly, symbolic AI struggles with real-world data’s unpredictability and variability. These challenges have led to the employment of statistical learning methodologies, like deep learning (Alom et al., 2018), which are more adept at managing noise and uncertain information through vector-valued representations.

**Sub-Symbolic Methods** The sub-symbolic framework, rooted in neural network paradigms, began with pioneering works such as the perceptron (McCulloch & Pitts, 1943), with the first hardware implementation quickly following (Rosenblatt, 1958). The foundational notion of distributed processing (Rumelhart et al., 1986) was later bolstered and further expanded by demonstrating that multilayer feedforward networks with a single hidden layer can serve as universal approximators for any Borel measurable function, given sufficient hidden units (Hornik et al., 1989). Fast-forward, contemporary frameworks achieve a significant leap with the introduction of the Transformer architecture (Vaswani et al., 2017), which underpins most of today’s LLMs. These LLMs demonstrate exceptional capabilities in in-context learning, a method popularized by the likes of GPT-3 (Brown et al., 2020), where models improve task performance through natural language instruction and examples provided directly in the input prompt. While in-context learning bypasses the need for explicit retraining, it demands meticulous prompt design to steer models towards desired behaviors.

**Neuro-Symbolic Methods** To overcome the limitations of each individual method, NeSy approaches meld the statistical inference strengths of deep neural architectures with the generalization and explainability of symbolic systems (Garcez et al., 2015; Besold et al., 2017; d’Avila Garcez et al., 2019; d’Avila Garcez & Lamb, 2020; Lamb et al., 2020; Hamilton et al., 2022; Yu et al., 2023). Some approaches focus on different strategies for integrating learning and reasoning processes (Yu et al., 2023; Fang et al., 2024). Firstly, *learning for reasoning* methods treat the learning aspect as an accelerator for reasoning, in which deep neural networks are employed to reduce the search space for symbolic systems (Silver et al., 2016; 2017b;a; Qu & Tang, 2019; Schrittwieser et al., 2020). Secondly, *reasoning for learning* views reasoning as a way to regularize learning, in which symbolic knowledge acts as a guiding constraint that oversees machine learning tasks (Hu et al., 2016; Xu et al., 2018). Thirdly, the *learning-reasoning* category enables a symbiotic relationship between learning and reasoning. Here, both elements interact and share information to boost

problem-solving capabilities (Donadello et al., 2017; Manhaeve et al., 2018; Mao et al., 2019; Ellis, 2023). This synergy further extends when considering graph-based methods, which closely align with the objectives of our proposed framework. Research in this area, such as CycleGT (Guo et al., 2020) and Paper2vec (Ganguly & Pudi, 2017) explores unsupervised techniques for bridging graph and text representations, GPTSwarm (Zhuge et al., 2024) explores graph optimizers to refine node-level prompts and edge optimization. Subsequently, graph embeddings, when utilized within symbolic frameworks, can enhance knowledge graph reasoning tasks (Zhang et al., 2021), or more generally, provide the bedrock for learning domain-invariant representations (Park et al., 2023).

Lastly, building upon the insights from Sun et al. (2022), the integration of NeSy techniques in scientific workflows promises significant acceleration in scientific discovery. While previous work has effectively identified opportunities and challenges, we have taken a more ambitious approach by developing a comprehensive framework from the ground up to facilitate a wide range of NeSy integrations.

**Large Language Models** In part, instruction-based fine-tuning of LLMs through reinforcement learning from human feedback (Ouyang et al., 2022; Li et al., 2023) or direct preference optimization (Rafailov et al., 2023) has shown promising results dealing with value misalignment issues (Bradley Knox & Stone, 2008; MacGlashan et al., 2017; Christiano et al., 2017; Ibarz et al., 2018; Goyal et al., 2022), unlocking new possibilities for chain of thoughts (Wei et al., 2022b), tree of thoughts (Yao et al., 2023a), and graph of thoughts interactions (Besta et al., 2023). However, recent research also highlights the limitations of LLMs in functional linguistic competence despite their proficiency in formal linguistic competence (Mahowald et al., 2023). Whereas formal linguistic competence encompasses the ability to understand and generate language, functional linguistic competence pertains to the application of language in real-world contexts, such as conveying sensory input or recalling information from memory. Examples of functional linguistic competence include implicatures (Ruis et al., 2022) and contextual language comprehension beyond the statistical manifestation of data distributions (Bransford & Johnson, 1972; Mikolov et al., 2013b). Consequently, operating LLMs through a purely inference-based approach confines their capabilities within their provided context window, severely limiting their horizon. This results in deficiencies for situational modeling, non-adaptability through contextual changes, and short-term problem-solving, amongst other capabilities. However, simply increasing the context length may not yield greater capabilities, as demonstrated by the observed U-shaped performance curve (Liu et al., 2023) where LLMs excel when utilizing information at the beginning or end of the input context, but struggle with information located in the middle, especially as context increases. These challenges are actively being researched, with novel approaches such as Hyena (Poli et al., 2023), RWKV (Bo, 2021), GateLoop (Katsch, 2023), Mamba (Gu & Dao, 2023) and xLSTM (Beck et al., 2024) surfacing. Meanwhile, the re-emergence of interest in retrieval-augmented generative approaches (Li et al., 2022a) offers an alternative by circumventing the autoregressive nature of the widely-utilized Transformer architecture (Vaswani et al., 2017), enabling context enrichment with lateral information.

**In-Context Learning** Recently, several in-context learning methodologies evolved to enable tool usage through LLMs (Schick et al., 2023), or refine the generative outcome of LLMs (Yang et al., 2023). This includes chain-of-thought (CoT) prompting, a method that conditions the model to reveal its step-by-step reasoning process (Wei et al., 2022b; Singhal et al., 2023). CoT prompting breaks down complex tasks into simpler, sequential steps, and helps with interpreting LLM’s output. Self-generated CoT, where models are encouraged to generate their own reasoning chains based on training examples, surpasses even expertly crafted CoT (Fernando et al., 2023). This observation echoes other reports that GPT-4 has an emergent self-improving capability through introspection, such as self-verification (Weng et al., 2023) or self-consistency (Wang et al., 2023b). Tree of Thoughts (ToT) enables LLMs to solve complex problems by exploring multiple reasoning paths through a search tree of coherent text units, demonstrating significant problem-solving enhancements in tasks requiring strategic planning and search (Yao et al., 2023a). Ensemble techniques further enhance the robustness and accuracy of model predictions by combining several strategies to establish a consensus (Nori et al., 2023).

### 3 PROBLEM DEFINITION

Conventional approaches employing foundation models, such as LLMs, are predominantly confined to single-step or few-step executions and primarily reliant on hand-crafted prompt instructions, often referred to as in-context learning. This restricted scope limits the utilization of different modalities, lacks verification, and exhibits limited tool proficiency. We posit that the use of NeSy engines as core computation units, realized through logic-based methodologies coupled with sub-symbolic foundation models, offers a more general, robust, and verifiable perspective. This approach has several advantages. Firstly, it enables the integration of pre-existing solutions (e.g. various classical algorithms), offloading computational complexity and bridging different modalities. Secondly, it allows sub-symbolic components to focus on decision-making (e.g. selecting the respective tool based on in-context classification). Thirdly, it provides

an *interpretable language-based control layer* for explainable, autonomous systems. In the following section, we elaborate on the key design principles underlying SymbolicAI and how we guide the generative processes of NeSy engines. For further technical details, see Appendix Section 5.

## 4 DESIGN PRINCIPLES

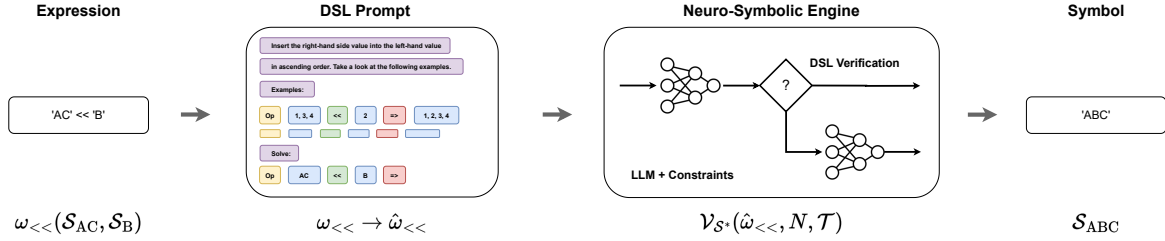
**Symbols and Expressions** As posited by [Newell & Simon \(1976\)](#), symbols are elemental carriers of meaning within a computational context<sup>3</sup>. These symbols define physical patterns capable of composing complex structures, and are central to the design and interpretation of logic and knowledge representations ([Augusto, 2022](#)). We define a symbol as the set  $\mathcal{S} = \bigcup_{n \geq 0} \mathbb{L}^n$  formed by concatenating characters from a finite character set  $\mathbb{L}$ , i.e. the vocabulary in an LLM setting, and with  $n$  representing the sequence length of the string. Thus, let the set of all possible symbols be defined as  $\Sigma$  and that  $\mathcal{S} \in \Sigma$ . We further introduce an operation  $\oplus$  that enables us to create expressions on any number of symbols from  $\Sigma$ , and when evaluated returns a new symbol in  $\Sigma$ . For any subset  $\{\mathcal{S}_1, \mathcal{S}_2, \dots, \mathcal{S}_m\} \subseteq \Sigma$ , an expression is defined as  $\omega : \bigoplus_{i=1}^m \mathcal{S}_i \rightarrow \mathcal{S}'$  from the set of all possible expressions  $\omega \in \Omega$ , where  $\mathcal{S}' \in \Sigma$ , and  $\bigoplus$  represents the placeholder operation of combining and transforming the symbols according to specific rules for  $m$  number of symbols. Such a specific rule for  $\bigoplus$  can define an arithmetic expression  $\bigoplus := +$  where two symbols are added, i.e.  $\omega := "1" + "two"$  which results in a new symbol "3" or "three". Thus, SymbolicAI is based on the concept that symbols, and the expressions they form, are reflections of the information inherent in a NeSy system, and serve as surrogate for the interaction between the NeSy system and the problem space. Moreover, we argue that *real patterns* ([Dennett, 1991](#)), recurring and identifiable structures that coherently and reliably emerge in the data beyond mere randomness or noise, can be effectively realized through symbols.

Furthermore, we utilize language as a tool for mapping complex concepts, leveraging its inherent semantics and abstractions to describe states and properties of a problem at hand. These mappings are universal, e.g. they may be utilized to define scene descriptions, long-horizon planning, acoustic properties, emotional states, physical conditions, etc. Therefore, language serves as a comprehensive, yet abstract framework to encapsulate meanings, and refer to it as the *convex hull of the knowledge of our society*. Subsequently, it is common to attribute existing physical objects with abstract concepts, as exemplified by our natural tendency to link tangible objects to colors and emotions, such as blending the color "red" with "heart", "warm", and "passion". This approach also anchors our work in the field of formal language theory, as we require a structured method to construct mappings from the world to language. Consequently, we use formal language structures, such as grammars, to systematically define our language-centric approach to problem-solving and the associated translation of real-world complexities into linguistic terms.

**Formal Languages** In formal language theory and linguistics, languages are structured following the Chomsky hierarchy, which classifies languages by the complexity of their grammatical structure ([Chomsky, 1956](#)). This hierarchy defines four types of grammars (Type-3 to Type-0) and separates formal languages by their grammatical complexity. A grammar in this context consists of terminal and non-terminal symbols, production rules, and a designated *start symbol*, enabling the generation of valid strings within a language.

We define a NeSy engine as a mapping  $\mathcal{V}_{\mathcal{S}^*} : \Omega \times N \times \mathcal{T} \rightarrow \Sigma$ , where  $N \subset \Sigma$  is a set of non-terminal symbols,  $\mathcal{T} \subset \Sigma$  is a set of terminal symbols and  $N \cap \mathcal{T} = \emptyset$ , and  $\mathcal{S}^* \in \Sigma$  is a starting symbol. We further formalize a grammar  $G = (N, \mathcal{T}, P, \mathcal{S}^*)$  with production rules defined as a  $P := \mathcal{V}_{\mathcal{S}^*}(\omega, N, \mathcal{T})$ . This grammar describes the generation of symbols through expressions  $\omega$ . For simplicity, we will drop the subscript of  $\mathcal{V}_{\mathcal{S}^*}$  and use it as  $\mathcal{V}$ . We identify LLMs as promising candidates for functioning as part of NeSy engines. In SymbolicAI, a symbol  $\mathcal{S}$  is augmented with conditional instructions and types derived from DSLs, custom defined or not (e.g. HTML, SQL, etc.), tailored for directing the LLMs. The key advantage of LLMs over previous systems lies in their ability to generalize across formal languages ([Wang et al., 2023a](#)) and knowledge systems. Although there is currently no universal consensus regarding the precise classification of natural language within the Chomsky hierarchy, our approach can be understood as employing a *situation-specific*, context-sensitive grammar, which enables the processing of instructions and analogies with a nuanced understanding of language. The intersection between formal and natural languages becomes evident when considering how language patterns, through prompts like "You are a helpful assistant...", elicit structured responses, indicating a potential underlying formal mechanism at play. This observation underlines the utility of such a grammar in our framework, where it serves as an explicit schema guiding the structure of examples for in-context learning. For instance, equating "3.1415..." with " $\pi$ " or "August 4, 1961" with "1961-08-04" in a given context demonstrates context-dependent interpretation of symbols. Such a system doesn't rigidly adhere to standard grammatical rules but instead adjusts and interprets based on the context, effectively creating a situation-specific gram-

<sup>3</sup> Our framework's name is derived from the foundational work of Newell and Simon.



**Figure 2:** Illustration for NeSy pipeline, showcasing conceptual usage of in-context learning methodologies, domain-specific language (DSL) structures, and the expression evaluations through a NeSy engine based on an LLM and constraint verification. The expression showcases the sorted insert operator  $\ll$  and how the information of the symbol B is included in the symbol AC. The violet placeholder in the *DSL Prompt* represents an instruction, such as "Insert the right-hand side value into the left-hand value in ascending order. Take a look at the following examples." The positions below represent task-specific few-shot examples. The DSL Prompt receives the expression  $\omega_{\ll}$  and maps it to  $\hat{\omega}_{\ll}$  that can be processed by the LLM-based NeSy function  $\mathcal{V}_{\mathcal{S}^*}$  and outputs a new symbol.

mar, capable of forming *Domain-Invariant Associations* through in-context learning. We further address this in a later paragraph.

**Function Composition** In SymbolicAI, we use function composition to construct complex hierarchies and behaviors from fundamental elements. Therefore, our framework enables modeling of interconnected processes, where the output of one function is used as input for another, thus creating a sequence of operations. Through function composition, we construct computational graphs, in which intermediate symbols represent the nodes or states within these graphs. Formally, function composition is denoted by  $\circ$ , where combining functions  $f$  and  $g$  yields a new function  $h = g \circ f$ , defined as  $h(x) = g(f(x))$ . For functions  $f : X \rightarrow Y$  and  $g : Y \rightarrow Z$ , their composition results in a function mapping elements from domain  $X$  to codomain  $Z$  through  $g(f(x))$ . Although traditionally the codomain of the inner function  $f$  aligns with the domain of the outer function  $g$ , SymbolicAI relaxes this constraint by allowing for any subset relationship between these domains and codomains, which is particularly beneficial for in-context learning. When using LLMs for NeSy production rules  $\mathcal{V}$ , we can derive a multi-step generative process by composing a computational graph as a sequence of zero- and few-shot function compositions:

$$\mathcal{V}(\omega_j, N, \mathcal{T}) = \mathcal{V}(\omega_{j-1}, \cdot) \circ \mathcal{V}(\omega_{j-2}, \cdot) \circ \dots \circ \mathcal{V}(\omega_0, \cdot), \quad (1)$$

where  $\omega_0$  is the initial instruction and  $j$  defines the index variable for a multi-step generative process. By leveraging functional in-context learning, where zero- and few-shot examples act as dynamic elements of the function's domain, SymbolicAI has the ability to interpret and respond to diverse input contexts. For instance, a function can classify a user request and select an appropriate interface (e.g. WolframAlpha) to process the request. The output modality may even vary based on the respective engine. This enables SymbolicAI to handle operations over multi-modal data that connects multi-step generative processes and establishes function composition as a central tenet in bridging multiple modalities and coordinating a variety of tasks.

**Domain-Invariant Associations** In-context learning enabled LLMs to become versatile task solvers by interpolating within the training distribution, to the extent that even potentially unseen tasks are addressable (Brown et al., 2020). We attribute this to associations formed within the input space and the capacity of Transformer architectures for defining domain-invariant feature sub-spaces. This phenomenon has parallels with few-shot learning approaches such as SubGD (Gauch et al., 2022), a method based on identifying and utilizing a low-dimensional subspace, learned from various tasks that effectively regularize the learning process. Since LLMs have been trained on different domains and tasks, which also include formulations of mathematical expressions, we posit that specific tokens, such as the equality sign, can be leveraged to associate meanings between different symbolic objects. Unlike domain-invariant representations that create invariant features across different learning tasks, our approach leverages the in-context generalization capability of LLMs to construct invariant symbolic associations that aim to preserve, manipulate and propagate situational context. We can use these properties to build operations that apply transformations on objects that are substitutes to the semantically aligned few-shot learning examples.

## 5 SYMBOLICAI FRAMEWORK

In this section, we discuss the specifics of the proposed SymbolicAI framework. For more details about the framework structure, see Appendix Section C. For installation and usage of our framework, see Appendix Section D. For more technical details and code snippets, see Appendix Section E.

**Types and Representations** Analogous to the type `object` in Python, the base type of SymbolicAI is a symbol represented by the base type `Symbol`. All other subtypes, such as `Expression`, represent their mathematical namesake and can be evaluated and simplified. These subtypes inherit from `Symbol` the base attributes, primitive operators, and helper methods.

Although SymbolicAI uses a language-centric design, modeling and manipulating every interaction into symbolic representations is not inherently efficient. Therefore, we establish mappings between symbolic and sub-symbolic representations for sensory inputs and non-discrete elements. Such mappings are typically realized through function approximation. This allows us to map between *modality*-to-language and language-to-*modality* use cases. Here, *modality* serves as a placeholder for various types such as text, image, video, audio, motion, etc. In turn, each `Symbol` object contains valued and vector-valued representations, obtained through `value` and `embedding` attributes. The latter represents a symbol’s current value, akin to embedding text and storing it as a PyTorch tensor (Paszke et al., 2019) or NumPy array (Harris et al., 2020). While for an LLM, the numerical tensors may lack inherent meaning, vector-valued representations play an important role when 1) composite symbols are combined into more complex expressions, and 2) these embedded tensors are updated through gradient-based optimization.

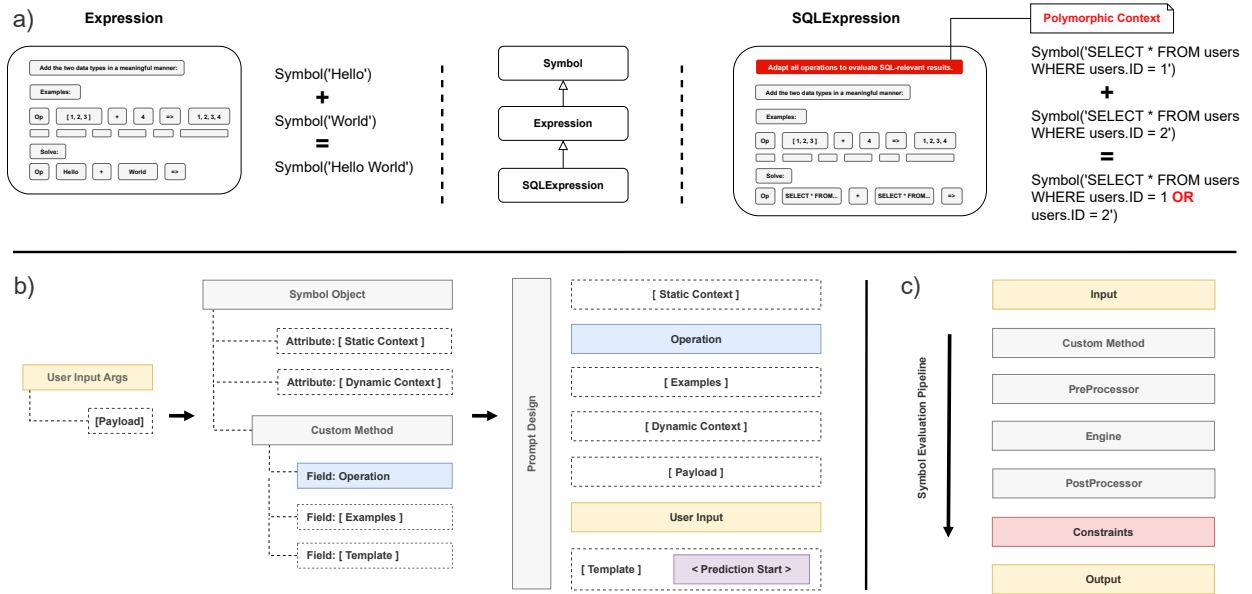
To enable the processing of symbols by LLMs, we assume that each `Symbol` object implements Python’s native string functionality, where the `__str__` method returns an interpretable string representation. Therefore, we can assert that any Python object is parsable by an LLM, however, the user must ensure a meaningful representation. For more details, see Appendix Section E.

**Polymorphic Context** Polymorphism is a central concept in programming language theory and prominently featured in SymbolicAI. Polymorphism refers to the ability of different objects to be accessed through the same interface, or of a single identifier to represent different types based on the context of execution. Providing a single interface for entities of different types allows operations to be performed in ways specific to their derived types. We therefore designed the `Symbol` object to contain a global context, which is composed of static and dynamic context parts, and enables this polymorphic behavior. The static context is class dependent and defined at design time. The dynamic context is runtime adaptable and can be changed to adhere to runtime specific logic and changes. Moreover, `Symbol` associated operations are resolved following polymorphic design before being evaluated by the NeSy engine. SymbolicAI’s engine implementation contains a `prepare` method to resolve and compile the engine specific representation by evaluating the `Symbol`-specific operations and context. For an example on polymorphic context see part a) in Figure 3.

**Operators and Methods** In SymbolicAI, operators are overloaded to facilitate transformations of `Symbol` objects. These operator primitives employ dynamic casting to ensure type compatibility. Consequently, `Symbol` objects can be easily manipulated through type specific attributions or symbolically evaluated by the NeSy engine. For example, a central operation for boolean logic is measuring equality between symbols. To evaluate the equality of symbols, we primarily adhere to the type specific implementation, because we prioritize strict comparisons over probabilistic evaluations. If the evaluation was unsuccessful, we then consider semantic equality through the NeSy engine. SymbolicAI leverages decorators for composing operators and custom class methods. For more details, see Appendix Section C.

Upon invoking an operator or method, the respective primitive function evaluates the symbol’s specific type and its respective attributes, and if necessary, resolves a nested decorated function that then uses the NeSy engine for evaluation. Should the evaluation fail, a predefined fallback implementation executes. Absent a fallback, or if both evaluations fail, an error state is raised. The processing of an operator or custom method involves a pipeline consisting of pre- and post-processing steps, as well as constraint enforcement. Constraints cover aspects like return types, value ranges, and structural integrity (e.g. JSON formatting through grammar-based verification). In Figure 3 b) we give an overview of the entire prompt composition based on the user input, the `Symbol` object structure, and in part c) the `Symbol` evaluation pipeline.

**Self-Referential Structures** SymbolicAI augments the generative process by enabling systems to introspect and modify their behavior dynamically. We leverage LLMs to execute tasks based on both natural and formal language instructions, adhering to the specified user objectives and with innate self-referential structures. We derive subtypes from `Expression` and enclose their functionalities in task-specific components, which we then expose again through



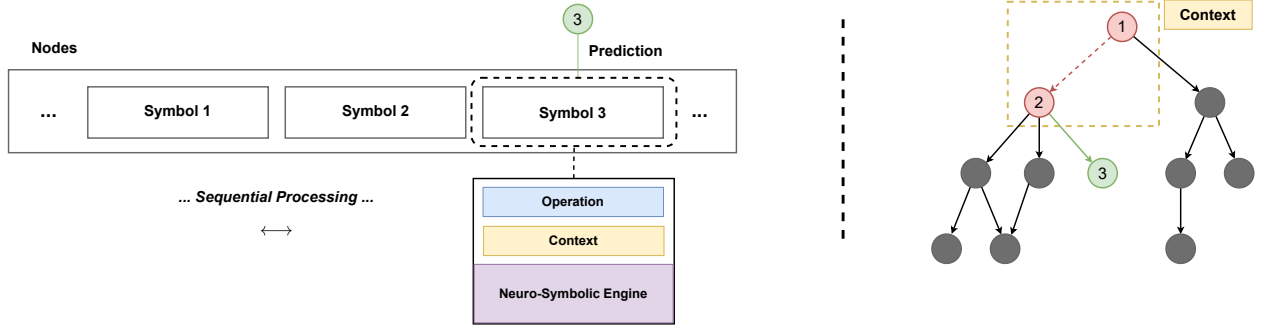
**Figure 3:** **a)** Illustration of polymorphic context on the example of a `SQLExpression` type for the `add`-operator. Without a polymorphic context a regular `Expression` evaluation concatenates two `Symbol` objects together. The polymorphic context in `SQLExpression` overwrites the base behavior such that two added SQL-expressions get semantically combined, not concatenated. **b)** Illustration of the translation of a `Symbol` object to a prompt statement to be processed by an LLM in the NeSy engine. The `User Input Args` can be attached with a `Payload` from previous executions and gets applied to the `Custom Method`. The user input with the polymorphic context of the `Symbol` Object attributes (`Static Context` and `Dynamic Context`) are translated to a prompt statement according to the schema of the `Prompt Design`. The fields `Operation`, `Examples` and `Template` mark operation description, DSL-based prompt examples and template structures respectively. These translations are processed according to `PreProcessor` and engine-specific formatting. **c)** Illustrates the evaluation pipeline from user input to output, with multiple translation processes before and after the `Engine` invocation. The `Input` gets passed to the `Custom Method` and reformatted according to a `PreProcessor` to adhere to DSL-specific structure. The engine then takes the output of the `PreProcessor` and composes the final prompt according to the engine-specific `Prompt Design` and resolves polymorphic context and auxiliary fields. The output of the `Engine` then can be restructured by a `PostProcessor` to match DSL-requirements of the desired `Output` and gets applied `Constraints` to verify the outcome.

templating and the model-driven design of the NeSy engine. This design choice allows a system to create and utilize its own sub-process definitions, analogous to concepts discussed in Schmidhuber (2007; 2009). Concretely, we utilize generalization properties of LLMs to interpret and formulate a set of operations that incorporate *self-instructions* (Wang et al., 2022). Consequently, the operations hold the flexibility to adapt to the context, and derive sub-processes that self-instruct LLMs to engage in situational modeling and context-sensitive problem-solving. Ultimately, this enables the construction of hierarchical computational graphs for self-referential *meta-reasoning* systems without the need to explicitly training a meta-learner (Kirsch & Schmidhuber, 2022). In Figure 4 we illustrate a step-wise evaluation of a contextual computational graph, in which the NeSy engine is processing conditioned on the current execution context and producing a next symbol prediction.

## 6 PERFORMANCE MEASURE

One of the challenges when creating multi-step generative processes with LLMs as part of NeSy engines relies on model evaluation and handling irrelevant predictions. The naïve assessment that measures only task succession would score all models to zero and render them as unusable. Even if models follow instructions and produce parts of the expected solution, we regularly observe that they — especially open-source models — append a continuation of task irrelevant predictions. Such predictions result in failure modes when applying conditions and validations, and halt any multi-step procedure. Our solution is an evaluation protocol that refines the performance measurement, allowing for more nuanced diagnostics and the possibility of continuing the evaluation despite intermediate failures. To derive our





**Figure 4:** We showcase a multi-step hierarchical computational graph, with each node in the graph represented by a symbol. The edges are relations between symbols. The left-hand side illustrates how a new node (Symbol 3) is obtained by evaluating an operation with its respective context on a NeSy engine. The right-hand side illustrates the context information window (yellow rectangle) and relationship of the resulting graph with its respective nodes.

quality measure, we borrow ideas from the utilization of the Fréchet distance for generative processes (Heusel et al., 2017).

We generate trajectories through a NeSy sequential process that creates a trajectory of distributions  $\mathbb{P}$  over multiple iterations of generative nodes. Each node in the process can be aligned to a reference distribution, which marks the desired behavior. To quantify the validity of the generated trajectories, we measure the total distance between the generated and reference data distribution along the path trajectory. We therefore adopt a cumulative measure capable of taking into account the entire generative trajectory. In theory, this process would entail calculating the path integral over the latent space representations for models, cumulating the Fréchet distances (Dowson & Landau, 1982) traversed along these trajectories:

$$\mathcal{D}(\mathbb{P}_{\text{gen}}, \mathbb{P}_{\text{ref}}) = \int_{t_0}^{t_f} d(\mathcal{N}(m_t, C_t), \mathcal{N}(m_{w,t}, C_{w,t})) dt \quad (2)$$

where  $\mathcal{D}(\mathbb{P}_{\text{gen}}, \mathbb{P}_{\text{ref}})$  denotes the integral of the Fréchet distances between two data distributions along the generative path trajectory from an initial time  $t_0$  to a final time  $t_f$ ,  $d(\mathcal{N}(m_t, C_t), \mathcal{N}(m_{w,t}, C_{w,t}))$  is the Fréchet distance calculated at each time  $t$  between the generated multivariate normal data distribution with mean  $m_t$  and covariance  $C_t$ , and the reference multivariate normal data distribution with mean  $m_{w,t}$  and covariance  $C_{w,t}$ . The resulting measure follows properties of normal distributions and is consistent with increasing disturbances.

However, this approach is computationally intractable for large-scale problems, and requires access to latent representations, which — especially in the context of LLMs — is not always given. For computational feasibility, we introduce an approximation that measures the embedding distances over the path trajectories through an auxiliary embedding model, based on prior work on distribution regression (Szabó et al., 2016). The embedding model maps the symbolic representations into a RKHS, such that we can apply a kernel mean embedding function to measure their respective distances (You et al., 2019; Dinu et al., 2023). We assess the distance through the mean embeddings w.r.t. to a kernel function  $K(\cdot, \cdot)$  of the samples  $\mathbf{e}_x^t \sim \nu_{\text{gen}}^t \in \mathbb{P}_{\text{gen}}$  and  $\mathbf{e}_y^t \sim \nu_{\text{ref}}^t \in \mathbb{P}_{\text{ref}}$  produced by the generated data distribution and a reference data distribution respectively. We denote by  $\mu_{\mathbf{e}_x^t}$ ,  $\mu_{\mathbf{e}_y^t}$  the mean embeddings associated to the respective samples, i.e.  $\mu_{\mathbf{e}_x^t}(z) = \frac{1}{n} \sum_{i=1}^n K(x_i^t, z)$  in case  $\mathbf{e}_x^t = (x_i^t)_{i=1}^n$  is a sample of size  $n$  of the respective mean embeddings. To compute the similarity between the embeddings of the generated and reference distributions, we evaluate the associated maximum mean discrepancy  $\text{MMD}^2(\mu_{\mathbf{e}_x^t}, \mu_{\mathbf{e}_y^t})$  (Gretton et al., 2012) and then, as before for the Fréchet distances, we integrate over  $t$ :

$$\tilde{\mathcal{D}}(\mathbb{P}_{\text{gen}}, \mathbb{P}_{\text{ref}}) = \int_{t_0}^{t_f} \text{MMD}^2(\mu_{\mathbf{e}_x^t}, \mu_{\mathbf{e}_y^t}) dt. \quad (3)$$

In empirical evaluations, however, we care about normalized values for ease of interpretation. We therefore analyze the properties of the MMD and derive a similarity score, which follows the same statistical principles as the MMD, and is bound between  $[0, 1]$ . We concluded that we can utilize only the MMD cross terms to evaluate the similarities. See Appendix Section B for more details. For our comparisons as referenced in Figure 6 we therefore denote the similarities rather than distances. We then come to the following formulation and refer to our empirical measure as the "Vector Embedding for Relational Trajectory Evaluation through Cross-similarity", or *VERTEX* score for short:

$$s(\mathbb{P}_{\text{gen}}, \mathbb{P}_{\text{ref}}) := \int_{t_0}^{t_f} \left[ \min(\max(0, \frac{1}{z} \widetilde{\text{MMD}}^2(\mu_{\mathbf{e}_x^t}, \mu_{\mathbf{e}_y^t}) - z_{\text{rand}}), 1) \right] dt. \quad (4)$$

We approximate the integral across time steps through Monte Carlo approximation. The introduced normalization constants denote the similarities to a random sequence  $z_{\text{rand}}$ , which functions as a baseline subtraction to recenter our results, and a given reference score to rescale w.r.t. to scores obtained from comparing related solutions  $z$ . Min-max scaling ensures the final measure is bounded between  $[0, 1]$ . This process reflects properties such as Hölder continuity that bounds the kernel function within certain limits. To compute the embeddings, we utilize the embedding model `all-mpnet-base-v2` (Song et al., 2020), due to its widespread availability, and its balance between speed and quality. As a similarity measure, we select a Gaussian kernel following our derivation from the Appendix Section B. In our implementations, we also explore other kernels, including preliminary experiments with cosine similarity. We also note that one can integrate Bernoulli distributed trials into our score, with 0 values representing failure modes and values of 1 being successes. Furthermore, if we relax our definition, we can integrate other similarity measures which are bound between  $[0, 1]$ , which then reflect on domain-specific attributions, i.e. including a similarity measure tailored towards capturing the nuances between two sub-structures of abstract syntax tree.

## 7 EVALUATION

We introduce a benchmark that evaluates multi-step generative processes as complex workflows. Our benchmark consists of five different evaluation categories, and uses the VERTEX score to measure the capabilities of an LLM to solve tasks from each category. The five categories of our benchmark are divided into three basic evaluations and two advanced categories that combine different basic capabilities. The three basic categories are (i) **associative prediction** which measures a models proficiency in understanding associations between symbols, (ii) **multi-modal binding** where we evaluate the capability to employ tools and operate on different modalities, and (iii) **program synthesis** for measuring a models proficiency in generating consistent code and executing subroutines. The two advanced benchmark categories are (iv) **logic**, for evaluating logic-based components and (v) **computational graphs** where complex workflows need to be processed, evaluating all aforementioned capabilities. For our evaluation we focus on the GPT family (Brown et al., 2020) of models, specifically GPT-3.5 Turbo (revision 1106) and GPT-4 Turbo (revision 1106) as they are the most proficient models to date; Gemini-Pro (Google, 2023) as the best performing model available through API from Google; LLaMA2-Chat 13B (Touvron et al., 2023), LLaMA3-Chat 8B and LLaMA3-Chat 70B from Meta represent open-source LLMs. Finally, Mistral 7B (Jiang et al., 2023) and Zephyr 7B (Tunstall et al., 2023) serve as baselines for revised and fine-tuned open-source models. The open-source models Mistral, Zephyr, and smaller LLaMA variants are estimated to have roughly equivalent parameter counts compared to GPT-3.5 Turbo and Gemini-Pro. All our experiments require a context size smaller or equal to 4096 to enable the comparisons among the in-context capabilities across model architectures. For the LLaMA models, we utilize the *chat* versions since they are specifically tuned to follow instructions.

**Associative Prediction** We evaluate a model’s proficiency to follow simple and complex instructions and associations with zero- and few-shot examples. Therefore, we evaluate the proficiency in applying our operators between `Symbol` types. We defined a total of 15 tasks involving in-context associations between two `Symbol` instances. SymbolicAI’s overloaded operators rely on predefined pseudo-grammars, as described in Section 4, that augment the operators with few-shot examples. For instance, the overloaded operator `+` utilized between two `Symbol` instances provides few-shot examples how to resolve additions with various data types. Consequently, we can now test if the models can solve the addition between `Symbol("two hundred and thirty four")` and `Symbol(7000)`. See Appendix Section F.1 for more details.

**Multi-modal Binding** We perform transformations between multiple modalities through language-based representations. Therefore, we need to evaluate the model’s proficiency in tool utilization, classification and routing of requests to relevant modules. We define a multi-modal `Expression` to detect the category of a task based on its content and to forward the task to the appropriate tool. The expression creates interfaces to tools like WolframAlpha for mathematical expressions, Selenium for website content scraping, SerpApi for search queries, and APILayer for optical character recognition. Each of the five tests aims to evaluate the appropriate handling of a specific type of input by the multi-modal `Expression` type, such as processing a website URL for scraping, interpreting a search engine query, testing if two vectors are linearly independent, comparing large numbers, and extracting text from an image. See Appendix Section F.2 for more details.

**Program Synthesis** We evaluate executable code with and without concepts from retrieval augmented generation, model-driven development, and experiment with self-generated instructions by creating self-referential expressions. We designed three separate tests related to program synthesis, where each task assesses the ability of the models to generate and execute code based on natural language instructions or provided templates:

- 1) The first task involves reading a LaTeX table template and data, then generating a function to populate the table with the given data.
- 2) The second task tests the automatic code generation for API calls by fetching data from a specified URL and extracting specific information from the retrieved content.
- 3) The third task evaluates the ability to construct a custom `Expression` that processes a `Symbol` through a specific `Function` component from the `SymbolicAI` package.

Each of the three tests follows a similar pattern, where the generated code is scored based on its similarity to valid references and normalized with random samples. See Appendix Section F.3 for more details.

**Logical Components** To evaluate the capabilities for logical reasoning of models, we condition them to create a sequence of expressions as self-contained components, and refer to higher-order logic for their assessment. Based on the underlying *type theory* originating from [Whitehead & Russell \(1925–1927\)](#), we evaluate a models’ capability to resolve statements in the form of *there exists  $x$  such that  $x$  fulfills  $y$* . Such quantifiers define the standard semantics of expressions, where their meaning is given by a semantic function. A semantic function maps a term from an abstract definition to a point in a domain, which is an interpretation of the term’s type and value. Therefore, these functions operate on types and values of expressions, and relations thereof. Subsequently, NeSy engines can formulate and evaluate at inference time logic-based instructions through Lisp, Prolog, or Mathematica ([McCarthy, 1959](#); [Colmerauer & Roussel, 1993](#); [Chen et al., 1993](#); [Inc., 2022](#)), or leverage solvers such as Z3 ([Moura & Bjørner, 2008](#)). Therefore, the result of a natural language statement when evaluated by a NeSy engine can be interpreted by any expert system which defines the corresponding semantic functions and process them either in a symbolic ([Feigenbaum et al., 1965](#); [Gamble et al., 1994](#)), differentiable ([Veličković & Blundell, 2021](#); [Ibarz et al., 2022](#)), or hybrid manner ([Kuncicky et al., 1991](#)).

We evaluate how proficient models are at interpreting custom DSLs and define expression statements. DSLs are designed to express logical relations and operations in a structured format, and supports human-readable and machine-interpretable formulations. The following example illustrates such relationships by translating a natural language statement into an expression statement, as follows:

Marvins has four paws and likes to meow when I pet its fur. Is Marvins a cat?

A DSL may enforce the usage of `HAS(·)`, `IS(·)`, etc. and may condition an LLM to produce the following expressions:

- `HasFourPaws( $x$ )`:  $x$  has four paws.
- `LikesToMeowWhenPetted( $x$ )`:  $x$  likes to meow when it is petted.
- `IsCat( $x$ )`:  $x$  is a cat.

These are then utilized to define the following logical expression:

$$\forall x (\text{HasFourPaws}(x) \wedge \text{LikesToMeowWhenPetted}(x) \Rightarrow \text{IsCat}(x)).$$

An automated theorem prover can now evaluate this statement for all values of  $x$  and assess the validity of the original query. Lastly, our evaluation uses symbolic mathematics to manipulate algebraic expressions. This involves defining symbols and performing operations like factorization, simplification, and algebraic manipulation. The symbols are placeholders for any value, enabling the definition of general expressions without specifying their values upfront.

We designed six tests to assess the logical capabilities of the candidate models and group them as follows. See Appendix Section F.4 for more details.

- 1) We utilize the Python library SymPy for symbolic mathematics to create the mathematical expression  $ax + bx - cx - ay - by + cy + d$ . The task for the model is then to factorize the expression and extract all unique symbols as a list.
- 2) Three tasks evaluate a models’ capability to resolve the logical operations AND, OR, and XOR. For instance, the test for logical AND combines the symbols `Symbol("The horn only sounds on Sundays")` and `Symbol("I hear the horn")` and compares the answer against the human-generated references *"The horn only sounds on Sundays and I hear the horn."* and *"Since I hear the horn it is Sunday."* Since there is a large number of possible solutions, there is high variability in the solution space. Each model might prefer a different solution.

3) For another task we use a custom `Expression` that defines a DSL syntax and semantic structure. We use this `Expression` to extract higher-order logic expressions from a natural language statement, namely the puzzle ‘Who is Jay’s brother?’<sup>4</sup>, that preserves the original relationships.

4) For the final task, we again use the puzzle ‘Who is Jay’s brother?’ to evaluate a models’ capability for complex conversions. We use the Z3 theorem prover (Moura & Bjørner, 2008) to solve the ‘Who is Jay’s brother’ puzzle conditioned on the Z3 solvers’ solution to Einsteins’ famous puzzle ‘Who owns the fish?’. The task involves an indirect translation from natural language to executable code by the Z3 solver; the solution to Einstein’s puzzle acts as a form of self-contained “documentation” for how the Z3 solver should be utilized. The test constructs a template, which includes the task instructions, puzzle statement, and reference to the Einstein’s puzzle solution. The models are then asked to analyze the given problem and solution format and create a Python function with Z3 syntax that can solve the ‘Who is Jay’s brother?’ puzzle. The dynamically generated code is executed within the test environment utilizing Python’s `exec` function. We check the access to the Z3 solver and run the generated `solve_puzzle` function supposed to contain the logic to solve the puzzle. Once executed, the assembled Z3 logical clauses are processed by the solver, which verifies that the set of constraints is satisfiable. If so, the model generated by the solver is queried for the puzzle’s solution and scored using our VERTEX score.

**Hierarchical Computational Graphs** We evaluate the capabilities of models to orchestrate a multi-step generative process and evaluate a set of tasks. Models need to direct sub-processes and associate computational results from and to `Symbol` nodes, and maintain relationships between these nodes, which we refer to as a computational graph as shown in Figure 5. In a computational graph, the VERTEX score compares the results produced by a generative model at each node against samples obtained from a reference distribution, usually modeled by sampling from multiple valid references. We also account for randomness through predefined random samples for normalizing the result. Our reference to *hierarchical* computational graphs stems from the fact that we operate on multiple levels. On a higher level of abstraction we are able to perform planning, sub-task scheduling, and define operational instructions. On a lower level of abstraction, we execute these plans based on the defined instructions and data, which can also span generative processes that produce new information.



**Figure 5:** We illustrate the hierarchical computational graph for the `Paper` expression. Each node represents an instance of an expression with distinct properties and behaviors, such as file sourcing, generative process, tool utilization, or transformation operation. The edges denote the reference relationships between expressions and indicate the flow of information. The blue highlighted nodes mark the main sequence nodes of expressions utilized to create parts of the paper draft, such as `Method` section, `RelatedWork` section, `Abstract` section, and so on. Each generative node is used for evaluating the VERTEX score. None-generative nodes such as search engine results are not evaluated, and we assume to obtain ground-truth values.

Given that the field is currently at an early stage in developing even sequential schedulers for LLM-based planning systems, our evaluations will be confined to sequential execution only. We introduce two tests designed to evaluate multi-step generative processes:

- 1) We simulate and evaluate the process of writing a research paper draft based on a predefined hierarchical computational graph that focuses on the content output of the computational graph rather than planning and scheduling functionality. See Appendix Section F.5 for more details.
- 2) We test the VERTEX Protocol as defined in Algorithm 1, which represents our general method for evaluating multi-step generative processes. We create a self-contained test scenario to illustrate an end-to-end evaluation and as a go-to reference for how our protocol can be deployed in a realistic environment. Our evaluation protocol is not only designed to analyze and score a series of instructions, but also to provide a structured basis for recording these

<sup>4</sup> Bob has two sons, John and Jay. Jay has one brother and father. The father has two sons. Jay’s brother has a brother and a father. Who is Jay’s brother?

processes. Furthermore, we note that our evaluation protocol is generally formulated, which allows the application of non-sequential planning and scheduling.

---

**Algorithm 1** VERTEX Protocol
 

---

**Require:** NeSy engine  $\mathcal{V}$  as an LLM, embedding engine  $\mathcal{E} : \Sigma \rightarrow \mathcal{H} \subset \mathbb{R}^d$ , symbols  $\{x_0, x^*, y^*\} \subset \Sigma$ , with  $x_0$  as the initial instruction,  $x^*$  as the payload resulted from executing  $\mathcal{V}$ ,  $y^*$  as the reference, and  $*$  acting as a placeholder for  $\mathcal{P}, \mathcal{T}, \mathcal{C}$ , capabilities  $\mathcal{C} = \{\mathcal{F}_1, \mathcal{F}_2, \mathcal{F}_3, \dots\}$ , where each  $\mathcal{F}_i$  represents a specific functional role within the system, plan  $\mathcal{P} \subset \Sigma$ , task  $\mathcal{T} \in \mathcal{P}$ , memory buffer  $\mathcal{M} \subset \Sigma$ , a scoring function  $\tilde{s} : \mathcal{H} \times \mathcal{H} \rightarrow [0, 1]$ , a scheduler  $\mathcal{Q}$ , an aggregator  $\mathcal{A}$ , and score variables  $\{s\} \in [0, 1]$ .

**Method:**

```

1:  $\mathcal{V}, \mathcal{E}, \mathcal{Q}, \mathcal{C}, y^{\mathcal{P}} \leftarrow \text{INIT}(\cdot)$                                 ▷ Initialize engines, scheduler, capabilities, expected plan.
2:  $\mathcal{M} \leftarrow \emptyset, \mathcal{A} \leftarrow \emptyset$                                 ▷ Initialize memory buffer and aggregator.
3:  $x^{\mathcal{P}} \leftarrow \text{GENERATEPLAN}(x_0, \mathcal{V})$                                 ▷  $\mathcal{V}$  generates plan based on initial instruction.
4:  $\text{EVALUATE}(x^{\mathcal{P}}, y^{\mathcal{P}}, \mathcal{E}, \mathcal{A}, \tilde{s})$                                 ▷ Embed, score, and aggregate plan similarity.
5:  $\mathcal{P}, \mathcal{M} \leftarrow \text{UNFOLDPLAN}(y^{\mathcal{P}}, \mathcal{M}, \mathcal{Q})$                         ▷  $\mathcal{Q}$  unfolds plan into actionable tasks and updates progression.
6: while  $\mathcal{P} \neq \emptyset$  do                                            ▷ Run until list of tasks is exhausted.
7:    $\mathcal{T}, y^{\mathcal{C}}, y^{\mathcal{T}} \leftarrow \text{SELECT}(\mathcal{M}, \mathcal{V})$                             ▷  $\mathcal{V}$  selects next task based on task progression.
8:    $\mathcal{F}_i \leftarrow \text{IDENTIFY}(\mathcal{T}, \mathcal{C}, \mathcal{V})$                                 ▷  $\mathcal{V}$  identifies task-related capability  $\mathcal{F}_i$ .
9:    $x^{\mathcal{C}}, x^{\mathcal{T}} \leftarrow \text{EXECUTE}(\mathcal{T}, \mathcal{F}_i, \mathcal{Q})$                         ▷  $\mathcal{Q}$  executes  $\mathcal{T}$  with capability  $\mathcal{F}_i$  and assign results  $x^{\mathcal{C}}, x^{\mathcal{T}}$ .
10:   $\text{EVALUATE}(x^{\mathcal{C}}, y^{\mathcal{C}}, x^{\mathcal{T}}, y^{\mathcal{T}}, \mathcal{E}, \mathcal{A}, \tilde{s})$                     ▷ Embed, score, and aggregate capability similarity.
11:   $\mathcal{P}, \mathcal{M} \leftarrow \text{UPDATE}(\mathcal{T}, \mathcal{P}, \mathcal{M}, \mathcal{Q})$                         ▷  $\mathcal{Q}$  updates plan and progression.
12: end while
13:  $s \leftarrow \text{FINALIZE}(\mathcal{A})$                                         ▷ Finalize aggregation of scores.
14: return  $s$                                                             ▷ Return aggregated score of plan execution.

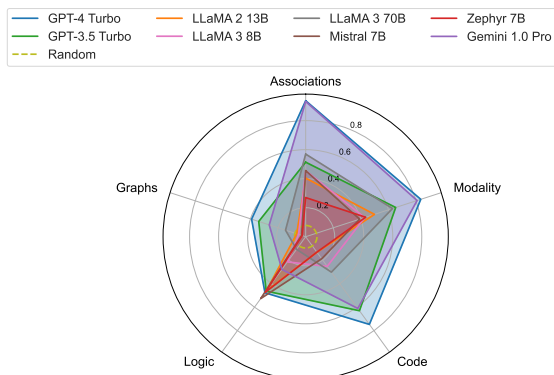
```

---

**Algorithm 1:** This algorithm defines the pseudocode of our VERTEX protocol with our respective VERTEX score as a scoring criteria. We start by initializing the NeSy engine  $\mathcal{V}$ , the embedding engine  $\mathcal{E}$ , the scheduler  $\mathcal{Q}$ , and a set of capabilities  $\mathcal{C}$ . The initial instruction  $x_0$  is utilized to generate a plan  $x^{\mathcal{P}}$  through  $\mathcal{V}$ . The plan and its expected outcome  $y^{\mathcal{P}}$  are embedded, and their similarity is scored according to our VERTEX score and aggregated. The plan is then unfolded into actionable tasks. Each task  $\mathcal{T}$  is selected and executed with the appropriate capability  $\mathcal{C}$ , resulting in the capability and task results  $x^{\mathcal{C}}, x^{\mathcal{T}}$ , and expected outcomes  $y^{\mathcal{C}}, y^{\mathcal{T}}$  updated in the memory buffer  $\mathcal{M}$ . The process continues, with each task's result being embedded, scored, and aggregated until the plan is complete. The final aggregated score  $s$  is returned, reflecting the overall effectiveness of the plan execution.

We start with a high-level workflow description which consists of a list of tasks and optionally their respective sub-tasks; we refer to this as the plan  $\mathcal{P}$ . To perform the experiment, we utilize an expected plan  $y^{\mathcal{P}}$  which was handcrafted for this evaluation. The expected plan is a queue of predefined tasks (in a particular order) that the system should follow to achieve the goal. The goal statement defines the end objective that the workflow aims to accomplish. We also have a set of plans similar to the expected plan, which are trajectories in the solution space, as well as the plan  $x^{\mathcal{P}}$  that the LLM generates utilizing the GENERATEPLAN call for a specific seed. We score the predicted plan against the expected plan and the trajectories, then we continue to the next phase in which we utilize the expected plan to execute the tasks. At each step, the LLM will receive in its context the goal, the tasks, the current progress, and a query asking for the next task to execute; we refer to this as the memory buffer  $\mathcal{M}$ . If the LLM is not able to predict the next task, it will return a failure, and the expected plan will be utilized to execute the next task. The LLM has access to a predefined set of capabilities  $\mathcal{C}$ , specifically WolframAlpha, SerpApi, Selenium, and the LLM itself, which also represents our self-referential structure. We keep executing tasks until the queue is exhausted, and at each step, we utilize the EVALUATE call to measure the performance of the LLM with our VERTEX score. The scheduler class  $\mathcal{Q}$  oversees the execution of the test workflow. It takes the setup configuration and orchestrates the linear execution of tasks, utilizing the expected plan as a reference. It maintains a pool of tasks to be executed and updates progress as tasks are completed. The UNFOLDPLAN call is a method of the scheduler class  $\mathcal{Q}$ . The method calls itself recursively until there are no tasks left. The SELECT call is responsible for determining which task to execute next from a pool of remaining tasks. It utilizes the LLM through self-reflection (Shinn et al., 2023) to choose the most suitable next task based on a template that gets progressively updated in the memory buffer  $\mathcal{M}$  by the UPDATE call. The IDENTIFY call uses self-reflection and similarity scoring to determine the best interface based on the task at hand, then passes the interface to the EXECUTE call to execute the task. Lastly, the test ends with the FINALIZE call, which provides an aggregated assessment of the model's ability to manage and execute the workflow.

In Figure 6 we conclude with our evaluation and compute the cumulative score for all described evaluation categories and in the next section we discuss how to interpret the results of our framework.



**Figure 6:** We evaluate GPT-4 Turbo, GPT-3.5 Turbo, Gemini-1.0 Pro, LLaMA2-Chat 13B, LLaMA3-Chat 8B, LLaMA3-Chat 70B, Mistral 7B and Zephyr 7B on five benchmark categories: 1) Associative Prediction (Association) 2) Multi-modal Binding (Modality) 3) Program Synthesis (Code) 4) Functional Logic Components (Logic) and 5) Hierarchical Computational Graphs (Graphs). We denote the VERTEX scores for each category as a normalized value between 0 and 1, where higher values are better. The VERTEX score is measured according to a reference baseline and normalized by random sequences to exclude noise and similarities among references distributions to rescale solutions. The shown scores are an average over all tests per category and across 8 different seeds per test.

Benchmarks	GPT-4 Turbo	GPT-3.5 Turbo	Gemini 1.0 Pro	LLaMA 2 13B	LLaMA 3 8B	LLaMA 3 70B	Mistral 7B	Zephyr 7B	Random
Associations	<b>0.94</b>	0.51	0.93	0.40	0.46	0.57	0.46	0.27	0.08
Modality	<b>0.83</b>	0.65	0.81	0.50	0.43	0.63	0.39	0.43	0.07
Code	<b>0.75</b>	0.63	0.61	0.13	0.25	0.30	0.19	0.13	0.00
Logic	0.48	0.46	0.28	0.46	0.21	0.11	<b>0.53</b>	0.47	0.00
Graphs	<b>0.39</b>	0.34	0.26	0.06	0.05	0.15	0.03	0.03	0.00
<b>Total</b>	<b>0.68</b>	0.52	0.58	0.31	0.28	0.35	0.32	0.27	0.03

## 8 DISCUSSION

In this section, we address the evaluation results, auxiliary findings and limitations of SymbolicAI and the future directions we are focusing on. Some of the limitations stem from the inherent constraints of current technologies and dependencies on third-party systems. Additionally, the nuanced complexities of working with generative models presents further challenges.

### 8.1 RESULTS

In Figure 6 we show the VERTEX score for all five evaluation categories on 8 different state-of-the-art models. We show the aggregated results per category, meaning the average score among all tests averaged per category and average across 8 different seeds per test. The VERTEX score is normalized between 0 and 1, where higher values are better. Our score is non-linear due to its nature of using non-linear kernels, and captures semantic, ordinal and relative structures among the data samples. However, since our score is highly dependent on the quality of the underlying embedding model, it may omit to capture fine-grained syntactic differences such as ‘Hello’ vs ‘hello’.

In our experiments, we have noticed that for associative predictions and multi-modal bindings, GPT-4 Turbo is on par with Gemini-1.0 Pro. Furthermore, there is still a large gap between open-source contestants such as LLaMA 3 even with 70B parameters compared to the closed-source alternatives from OpenAI and Google. For the rest of the experiments, we see that GPT-4 almost always outperforms all other models, except for the functional logic components category. Here, we analyzed results and found that the larger models sometimes take shortcuts by automatically returning the solution and answering that the task instructions are too complex for such a straight-forward puzzle query. However, we would rather state in general that for logic-based, planning and scheduling tasks all models act unreliably, even if slight performance differences between the models are seen in the plot. We believe this is in part due to lack of training data specifically for workflows, planning and scheduling tasks, and to imprecision in generating reliably structured output formats, such as custom DSLs or other custom in-context instructed formats. This also stems from their instruction fine-tuning, since most models are chat-based models and offer verbose responses which need to be suppressed or post-processed.

We see similar performance between GPT-3.5 Turbo and LLaMA 3 70B except for the logical and graphs evaluations. We found that LLaMA 3 70B has a tendency to ask questions back if it does not understand the request instead of following the specified instructions provided. We assume this also stems from the chat-based instruction fine-tuning. Zephyr 7B and Mistral 7B have shown on par capabilities in functional logic components with larger models, however

fail in program synthesis and hierarchical computational graphs experiments. We observe that they perform well when resolving the overloaded logic operators such as OR, AND and XOR, and show decent performance for text generation, but fail to resolve more complex instructions.

## 8.2 LIMITATIONS

**Framework** Since the framework interfaces with many tools and API services, it requires a substantial engineering feat to integrate all available functionalities and keep the API-based services up-to-date. For us this means, that although we support a variety of tools and frameworks like Selenium, WolframAlpha, or Z3, we only scratch the surface of these tools. Moreover, the utilization of grammar-based constraints validations is still experimental and limited in functionality for specific formats such as JSON and HTML. Finally, we encounter also challenges related to engineering parallelization and multiprocessing of prompts, since the concurrent execution is non-trivial, especially with intricacies of Python process management.

**Embedding Measure** Our empirical measure is limited by the expressiveness of the embedding model and how well it captures the nuances in similarities between two representations. Furthermore, the obtained similarity scores are highly non-linear and difficult to interpret. For instance, two representations may address the same topic, such as the problem description and its respective solution, however, when measuring their similarity we obtain similarity scores of  $\sim 70\%$ . We normalize this by subtracting an inherent baseline and randomness effect, however, to ensure a more holistic and robust measurement we would need a significantly larger amount of baselines and experiments. Since we were very limited in the availability of development resources, and some presented models are only addressable through costly API walls. We are actively seeking sponsors to scale our solution and offer a more compelling benchmark suite in the future.

**Model Capabilities** An obvious limitation revolves around the fixed context window size of the underlying language models. Despite the expansion of the context window in newer models such as GPT-4, the finite context still restricts the amount of data that can be processed in a single pass. All information outside the context needs to be added through information retrieval approaches, which come with their own challenges and limitations (Gao et al., 2023). This leads to side effects, including hallucination, given the model does not contain the necessary information to answer the prompted instruction, which makes it difficult to maintain long-term statefulness for complex reasoning tasks and computational graphs.

**Error Handling** The complexity of error handling when evaluating complex expressions through function compositionality, especially between multiple modalities and different solvers, is another notable challenge. While SymbolicAI introduces mechanisms for error analysis and automated correction, these approaches are not infallible. They are often limited by the quality and expressiveness of the models, and the model’s capacity to understand deeply nested logical constructs. We also note that for our evaluations, we disabled any remedy protocol, such as truncating prompts or retry schema.

**Generalization** This research is also limited by current LLM’s capacity for reasoning and generalization. Although progress has been made, models are still prone to hallucinations and reasoning errors, especially when dealing with abstract, novel, or highly complex problem statements (Marcus, 2020). Furthermore, our framework’s reliance on the model’s ability to grasp the semantics of operations can be influenced by the training data and the model’s innate biases and conceptual understanding (Mahowald et al., 2023). We also point out that the initial development of SymbolicAI started with the GPT family of models, and we may encounter innate biases in prompt design and expressiveness when utilizing other reference models. However, we also point out that prompt engineering instruction-based statements is not a reliable direction for improvement. We instead advocate for enhancing the resilience of models through fault tolerance, focusing on their ability to better follow semantic instructions, not syntactic idiosyncrasies. Another concern is how to assess the disentanglement of evaluations of models on downstream tasks, to avoid testing on training samples, especially for closed-source solutions like GPT.

**Interpretability and Transparency** Finally, the issue of explainability and transparency in AI systems remains challenging. While SymbolicAI makes steps towards making computational processes more explicit and explainable through symbolic manipulations, understanding the internal logic and decision-making of LLMs remains an open problem. This can hinder trust and adoption in sensitive applications where interpretability of predictions is important.

### 8.3 FUTURE WORK

The goal for Algorithm 1 is to be utilized by an advanced learning agent. This agent, employing reinforcement learning methodologies (Ouyang et al., 2022; Li et al., 2023; Rafailov et al., 2023), could leverage our evaluation measure in existing benchmarks (Milani et al., 2020; Swazinna et al., 2022; Schweighofer et al., 2022) as a means to obtain reward signals to address a central problem in reinforcement learning, namely credit assignment (Sutton, 1984; Arjona-Medina et al., 2019; Holzleitner et al., 2020; Patil et al., 2020; Widrich et al., 2021; Dinu et al., 2022). Over time, it aims to develop the ability to autonomously generate its own plans, efficiently schedule tasks and subtasks, and carefully select the most suitable tools for each task. Our protocol lays the groundwork for this agent to learn and expand its base set of capabilities (Amaro et al., 2023), moving towards more sophisticated, self-referential orchestration of multi-step tasks. We’ve already noticed that research is shifting towards this type of methodology (Yuan et al., 2024). Furthermore, in Section 7 we’ve only considered a sequential scheduler. However, our objective is to ultimately assess a non-sequential task execution model, allowing for dynamic insertion and out-of-sequence task execution. In addition, we are interested into exploring similarities of our work with *Generative Flow Networks* (Bengio et al., 2021a;b; Lahlou et al., 2023). Lastly, we also discuss limitations in Appendix Section 8.2 with further opportunities for future improvements.

## 9 CONCLUSION

In this work, we introduced SymbolicAI, a framework that unifies generative models with an array of solvers, blending the strengths of symbolic and sub-symbolic AI paradigms within a cohesive NeSy framework. SymbolicAI equips researchers and practitioners with a comprehensive toolkit to develop contextualized and explainable NeSy AI systems capable of addressing diverse challenges effectively. We also introduce a quality measure and a benchmark for comparing and evaluating a wide range of computational tasks. SymbolicAI provides a basis for further research in advanced program synthesis, hierarchical computational graphs, the development of self-referential systems, and the integration of probabilistic methods with AI design for creating autonomous agents.

## ACKNOWLEDGEMENT

The ELLIS Unit Linz, the LIT AI Lab, the Institute for Machine Learning, are supported by the Federal State Upper Austria. We thank the projects Medical Cognitive Computing Center (MC3), INCONTROL-RL (FFG-881064), PRIMAL (FFG-873979), S3AI (FFG-872172), DL for GranularFlow (FFG-871302), EPILEPSIA (FFG-892171), AIRI FG 9-N (FWF-36284, FWF-36235), AI4GreenHeatingGrids (FFG- 899943), INTEGRATE (FFG-892418), ELISE (H2020-ICT-2019-3 ID: 951847), Stars4Waters (HORIZON-CL6-2021-CLIMATE-01-01). We thank Audi.JKU Deep Learning Center, TGW LOGISTICS GROUP GMBH, Silicon Austria Labs (SAL), FILL Gesellschaft mbH, Anyline GmbH, Google, ZF Friedrichshafen AG, Robert Bosch GmbH, UCB Biopharma SRL, Merck Healthcare KGaA, Verbund AG, GLS (Univ. Waterloo), Software Competence Center Hagenberg GmbH, Borealis AG, TÜV Austria, Frauscher Sonsonic, TRUMPF, the NVIDIA Corporation and Atlas.

We extend our appreciation to Andreas Windisch and Clemens Wasner of AI Austria for their unwavering support. Their valuable feedback, connections, and facilitation of introductions within their expansive network have been instrumental to the progress of ExtensityAI.

Our gratitude also goes to Sergei Pereverzyev, whose enlightened guidance and thoughtful ideas have been a beacon for our research endeavors. Our thanks are equally extended to Gary Marcus, whose stimulating discussions sparked numerous innovative ideas incorporated into our framework.

We are equally grateful to Markus Hofmarcher, a friend and colleague whose informed counsel and stimulating discussions have significantly sharpened various facets of our study. Additionally, our thanks are due to Fabian Paischer and Kajetan Schweighofer, whose preliminary work and assistance have been of enormous benefit.

We are also grateful to our friends John Chong Min Tan and Tim Scarfe, whose communities have been a hub for exhilarating discussions. Their online presence and engagement have enriched the AI research landscape and broadened our perspectives.

Moreover, we wish to honor the memories of the cherished family members we lost in 2023. Their influence in our lives extended beyond personal bonds, and the principles they instilled in us continue to shape our journey. It is with great respect and affection that we acknowledge the indelible impact they have made, enabling us to persist in our scientific pursuits with determination and integrity.



## REFERENCES

- M. Z. Alom, T. M. Taha, C. Yakopcic, S. Westberg, P. Sidike, M. S. Nasrin, B. C. Van Esesn, A. A. S. Awwal, and V. K. Asari. The history began from alexnet: A comprehensive survey on deep learning approaches. *arXiv preprint arXiv:1803.01164*, 2018.
- R. E. Amaro, J.-Y. Chen, J. M. Duarte, T. E. Hutton, C. Irving, M. C. Kandes, A. Majumdar, D. Y. Mishin, M. H. Nguyen, P. Rodriguez, F. Silva, R. S. Sinkovits, S. M. Strande, M. Tatineni, L. S. Tran, and N. Wolter. Voyager – an innovative computational resource for artificial intelligence & machine learning applications in science and engineering. In *Practice and Experience in Advanced Research Computing, PEARC '23*, pp. 278–282, New York, NY, USA, 2023. Association for Computing Machinery. ISBN 9781450399852. doi: 10.1145/3569951.3597597.
- Peter Anderson, Basura Fernando, Mark Johnson, and Stephen Gould. SPICE: semantic propositional image caption evaluation. *CoRR*, abs/1607.08822, 2016. URL <http://arxiv.org/abs/1607.08822>.
- J. Andreas. Language models as agent models. *CoRR*, abs/2212.01681, 2022. doi: 10.48550/arXiv.2212.01681.
- J. A. Arjona-Medina, M. Gillhofer, M. Widrich, T. Unterthiner, J. Brandstetter, and S. Hochreiter. RUDDER: return decomposition for delayed rewards. In *Advances in Neural Information Processing Systems 32*, pp. 13566–13577, 2019.
- M. Assran, Q. Duval, I. Misra, P. Bojanowski, P. Vincent, M. Rabbat, Y. LeCun, and N. Ballas. Self-supervised learning from images with a joint-embedding predictive architecture. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 15619–15629, 2023.
- L. M. Augusto. *Computational Logic. Vol. 1: Classical Deductive Computing with Classical Logic*. College Publications, London, 2 edition, 2022.
- F. Badita. *1337 Use Cases for ChatGPT & other Chatbots in the AI-Driven Era*. Google Docs, 2022.
- D.M. Beazley. *Python Essential Reference*. Developer’s library : essential references for programming professionals. Addison-Wesley, 2009. ISBN 9780672329784. URL <https://books.google.ro/books?id=Chr1ND1UcI8C>.
- M. Beck, K. Pöppel, M. Spanring, A. Auer, O. Prudnikova, M. Kopp, G. Klambauer, J. Brandstetter, and S. Hochreiter. xlstm: Extended long short-term memory, 2024.
- E. Bengio, M. Jain, M. Korablyov, D. Precup, and Y. Bengio. Flow network based generative models for non-iterative diverse candidate generation. *Advances in Neural Information Processing Systems*, 34:27381–27394, 2021a.
- Y. Bengio, T. Deleu, E. J. Hu, S. Lahlou, M. Tiwari, and E. Bengio. Gflownet foundations. *arXiv preprint arXiv:2111.09266*, 2021b.
- T. R. Besold, A. d. Garcez, S. Bader, H. Bowman, P. Domingos, P. Hitzler, K.-U. Kuehnberger, L. C. Lamb, D. Lowd, P. M. V. Lima, L. de Penning, G. Pinkas, H. Poon, and G. Zaverucha. Neural-symbolic learning and reasoning: A survey and interpretation, 2017.
- M. Besta, N. Blach, A. Kubicek, R. Gerstenberger, L. Gianinazzi, J. Gajda, T. Lehmann, M. Podstawski, H. Niewiadomski, P. Nyczyk, and T. Hoefler. Graph of thoughts: Solving elaborate problems with large language models. *arXiv preprint arXiv:2308.09687*, 2023.
- S. Biderman, H. Schoelkopf, Q. Anthony, H. Bradley, K. O’Brien, E. Hallahan, M. Aflah Khan, S. Purohit, S. Prashanth, E. Raff, A. Skowron, L. Sutawika, and O. van der Wal. Pythia: A suite for analyzing large language models across training and scaling, 2023.
- PENG Bo. Blinkdl/rwkv-lm: 0.01. Technical report, Zenodo, August 2021. URL <https://doi.org/10.5281/zenodo.5196577>.
- W. Bradley Knox and Peter Stone. TAMER: Training an Agent Manually via Evaluative Reinforcement. In *2008 7th IEEE International Conference on Development and Learning*, pp. 292–297, Monterey, CA, August 2008. IEEE. ISBN 978-1-4244-2661-4. doi: 10.1109/DEVLRN.2008.4640845.
- J. D. Bransford and M. K. Johnson. Contextual prerequisites for understanding: Some investigations of comprehension and recall. *Journal of Verbal Learning and Verbal Behavior*, 11(6):717–726, 1972. ISSN 0022-5371.

- T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, T. Henighan, R. Child, A. Ramesh, D. Ziegler, J. Wu, C. Winter, C. Hesse, M. Chen, E. Sigler, M. Litwin, S. Gray, B. Chess, J. Clark, C. Berner, S. McCandlish, A. Radford, I. Sutskever, and D. Amodei. Language models are few-shot learners. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 1877–1901. Curran Associates, Inc., 2020.
- H. Chase. LangChain. Technical report, LangChain, 01 2023. URL <https://github.com/hwchase17/langchain>.
- W. Chen, M. Kifer, and D. S. Warren. Hilog: A foundation for higher-order logic programming. *The Journal of Logic Programming*, 15(3):187–230, 1993. ISSN 0743-1066.
- F. Chollet. On the measure of intelligence, 2019.
- N. Chomsky. Three models for the description of language. *IRE Transactions on Information Theory*, 2(3):113–124, 1956. doi: 10.1109/TIT.1956.1056813.
- P. F. Christiano, J. Leike, T. Brown, M. Martic, S. Legg, and D. Amodei. Deep reinforcement learning from human preferences. *Advances in neural information processing systems*, 30, 2017.
- A. Colmerauer and P. Roussel. The birth of Prolog. In *HOPL-II*, 1993.
- A. d’Avila Garcez and L. C. Lamb. Neurosymbolic ai: The 3rd wave. *arXiv preprint arXiv:2012.05876*, 2020.
- A. d’Avila Garcez, M. Gori, L. C. Lamb, L. Serafini, M. Spranger, and S. N. Tran. Neural-symbolic computing: An effective methodology for principled integration of machine learning and reasoning. *Journal of Applied Logic*, 2019.
- A. Dawid and Y. LeCun. Introduction to latent variable energy-based models: A path towards autonomous machine intelligence. *arXiv preprint arXiv:2306.02572*, 2023.
- J. Degraeve. Building A Virtual Machine inside ChatGPT. Technical report, Engraved, 11 2022. URL <https://www.engraved.blog/building-a-virtual-machine-inside/>.
- F. Deniz, A. O. Nunez-Elizalde, A. G. Huth, and J. L. Gallant. The representation of semantic information across human cerebral cortex during listening versus reading is invariant to stimulus modality. *Journal of Neuroscience*, 39(39):7722–7736, 2019. ISSN 0270-6474. doi: 10.1523/JNEUROSCI.0675-19.2019.
- D. C. Dennett. Real patterns. *Journal of Philosophy*, 88(1):27–51, 1991. doi: 10.2307/2027085.
- M. Dilhara, A. Ketkar, and D. Dig. Understanding software-2.0: A study of machine learning library usage and evolution. *ACM Transactions on Software Engineering and Methodology (TOSEM)*, 30(4):55:1–55:42, jul 2021. ISSN 1049-331X. doi: 10.1145/3453478.
- M.-C. Dinu, M. Hofmarcher, V. P. Patil, M. Dorfer, P. M. Blies, J. Brandstetter, J. A. Arjona-Medina, and S. Hochreiter. Xai and strategy extraction via reward redistribution. In A. Holzinger, R. Goebel, R. Fong, T. Moon, K.-R. Müller, and W. Samek (eds.), *xxAI - Beyond Explainable AI: International Workshop, Held in Conjunction with ICML 2020, July 18, 2020, Vienna, Austria, Revised and Extended Papers*, pp. 177–205, Cham, 2022. Springer International Publishing. ISBN 978-3-031-04083-2. doi: 10.1007/978-3-031-04083-2\_10.
- M.-C. Dinu, M. Holzleitner, M. Beck, H. D. Nguyen, A. Huber, H. Eghbal-zadeh, B. A. Moser, S. V. Pereverzyev, S. Hochreiter, and W. Zellinger. Addressing parameter choice issues in unsupervised domain adaptation by aggregation. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023.
- I. Donadello, L. Serafini, and A. d’Avila Garcez. Logic tensor networks for semantic image interpretation. In *Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence, IJCAI-17*, pp. 1596–1602, 2017.
- D. C. Dowson and B. V. Landau. The fréchet distance between multivariate normal distributions. *Journal of Multivariate Analysis*, 12(3):450–455, 1982. doi: [https://doi.org/10.1016/0047-259X\(82\)90077-X](https://doi.org/10.1016/0047-259X(82)90077-X).
- Kevin Ellis. Human-like few-shot learning via bayesian reasoning over natural language. *arXiv preprint arXiv:2306.02797*, 2023.

- M. Fang, S. Deng, Y. Zhang, Z. Shi, L. Chen, M. Pechenizkiy, and J. Wang. Large language models are neurosymbolic reasoners. *arXiv preprint arXiv:2401.09334*, 2024.
- E. Fedorenko, P.-J. Hsieh, A. Nieto-Castanon, S. Whitfield-Gabrieli, and N. Kanwisher. New method for fMRI investigations of language: Defining rois functionally in individual subjects. *Journal of neurophysiology*, 104:1177–94, 08 2010. doi: 10.1152/jn.00032.2010.
- E. Feigenbaum, B. G. Buchanan, J. Lederberg, Carl Djerassi, and et al. Dendral, 1965.
- C. Fernando, D. Banarse, H. Michalewski, S. Osindero, and T. Rocktäschel. Promptbreeder: Self-referential self-improvement via prompt evolution. *arXiv preprint arXiv:2309.16797*, 2023.
- R. F. Gamble, G.-C. Roman, H. C. Cunningham, and W. E. Ball. Applying formal verification methods to rule-based programs. *Int. J. Expert Syst.*, 7(3):203–237, sep 1994. ISSN 0894-9077.
- S. Ganguly and V. Pudi. Paper2vec: Combining graph and text information for scientific paper representation. In Joemon Jose et al. (eds.), *Advances in Information Retrieval*, volume 10193 of *Lecture Notes in Computer Science*. Springer, Cham, 2017. ISBN 978-3-319-56607-8. doi: 10.1007/978-3-319-56608-5\_30.
- Y. Gao, Y. Xiong, X. Gao, K. Jia, J. Pan, Y. Bi, Y. Dai, J. Sun, and H. Wang. Retrieval-augmented generation for large language models: A survey. *arXiv preprint arXiv:2312.10997*, 2023.
- A. Garcez, T. Besold, L. De Raedt, P. Földiák, P. Hitzler, T. Icard, K. Kühnberger, L. Lamb, R. Miikkulainen, and D. Silver. Neural-symbolic learning and reasoning: Contributions and challenges. In *AAAI Conference*, 2015.
- M. Gauch, M. Beck, T. Adler, D. Kotsur, S. Fiel, H. Eghbal-zadeh, J. Brandstetter, J. Kofler, M. Holzleitner, W. Zellinger, D. Klotz, S. Hochreiter, and S. Lehner. Few-Shot Learning by Dimensionality Reduction in Gradient Space. *arXiv preprint arXiv:2206.03483*, 2022.
- X. Geng, A. Gudibande, H. Liu, E. Wallace, P. Abbeel, S. Levine, and D. Song. Koala: A dialogue model for academic research. Blog post, April 2023. URL <https://bair.berkeley.edu/blog/2023/04/03/koala/>.
- Google. Gemini: A family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
- A. Goyal, A. Friesen, A. Banino, T. Weber, N. R. Ke, A. P. Badia, A. Guez, M. Mirza, P. C. Humphreys, K. Konyushova, M. Valko, S. Osindero, T. Lillicrap, N. Heess, and C. Blundell. Retrieval-augmented reinforcement learning. In K. Chaudhuri, S. Jegelka, L. Song, C. Szepesvari, G. Niu, and S. Sabato (eds.), *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pp. 7740–7765. PMLR, 17–23 Jul 2022.
- A. Gretton, K. M. Borgwardt, M. J. Rasch, B. Schölkopf, and A. Smola. A kernel two-sample test. *Journal of Machine Learning Research*, 13(25):723–773, 2012.
- A. Gu and T. Dao. Mamba: Linear-time sequence modeling with selective state spaces. *arXiv preprint arXiv:2312.00752*, 2023.
- Q. Guo, Z. Jin, X. Qiu, W. Zhang, D. Wipf, and Z. Zhang. CycleGT: Unsupervised graph-to-text and text-to-graph generation via cycle training. *arXiv preprint arXiv:2006.04702*, 2020.
- K. Hamilton, A. Nayak, B. Božić, and L. Longo. Is neuro-symbolic AI meeting its promises in natural language processing? a structured review. *Semantic Web*, pp. 1–42, nov 2022. doi: 10.3233/sw-223228.
- C. R. Harris, K. J. Millman, S. J. van der Walt, R. Gommers, P. Virtanen, D. Cournapeau, E. Wieser, J. Taylor, S. Berg, N. J. Smith, R. Kern, M. Picus, S. Hoyer, M. H. van Kerkwijk, M. Brett, A. Haldane, J. Fernández del Río, M. Wiebe, P. Peterson, P. Gérard-Marchant, K. Sheppard, T. Reddy, W. Weckesser, H. Abbasi, C. Gohlke, and T. E. Oliphant. Array programming with NumPy. *Nature*, 585(7825):357–362, 2020. doi: 10.1038/s41586-020-2649-2.
- M. Heusel, H. Ramsauer, T. Unterthiner, B. Nessler, and S. Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, NIPS’17, pp. 6629–6640, Red Hook, NY, USA, 2017. Curran Associates Inc. ISBN 9781510860964.
- S. Hochreiter. Toward a broad AI. *Commun. ACM*, 65(4):56–57, mar 2022. ISSN 0001-0782.
- S. Hochreiter and J. Schmidhuber. Flat minima. *Neural Comput.*, 9(1):1–42, 1997.

- M. Holzleitner, L. Gruber, J. A. Arjona-Medina, J. Brandstetter, and S. Hochreiter. Convergence proof for actor-critic methods applied to PPO and RUDDER. *arXiv preprint arXiv:2012.01399*, 2020.
- K. Hornik, M. Tinchcombe, and H. White. Multilayer feedforward networks are universal approximators. *Neural Networks*, 2:359–366, 1989. doi: 10.1016/0893-6080(89)90020-8.
- J. Hu, H. Small, H. Kean, A. Takahashi, L. Zekelman, D. Kleinman, E. Ryan, A. Nieto-Castañón, V. Ferreira, and E. Fedorenko. Precision fMRI reveals that the language-selective network supports both phrase-structure building and lexical access during language production. *bioRxiv*, 2022. doi: 10.1101/2021.09.10.459596.
- Z. Hu, X. Ma, Z. Liu, E. Hovy, and E. Xing. Harnessing deep neural networks with logic rules. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 2410–2420, Berlin, Germany, August 2016. Association for Computational Linguistics.
- A. G. Huth, W. A. de Heer, T. L. Griffiths, F. E. Theunissen, and J. L. Gallant. Natural speech reveals the semantic maps that tile human cerebral cortex. *Nature*, 532(7600):453–458, 2016. doi: 10.1038/nature17637.
- B. Ibarz, J. Leike, T. Pohlen, G. Irving, S. Legg, and D. Amodei. Reward learning from human preferences and demonstrations in atari. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018.
- B. Ibarz, V. Kurin, G. Papamakarios, K. Nikiforou, M. Abbana Bennani, R. Csordás, A. Dudzik, M. Bošnjak, A. Vitvitskiy, Y. Rubanova, A. Deac, B. Bevilacqua, Y. Ganin, C. Blundell, and P. Velivčković. A generalist neural algorithmic learner. In *LOG IN*, 2022.
- Wolfram Research, Inc. Mathematica, Version 13.2, 2022. URL <https://www.wolfram.com/mathematica>. Champaign, IL.
- G. Indiveri, B. Linares-Barranco, T. Hamilton, A. van Schaik, R. Etienne-Cummings, T. Delbruck, S. Liu, P. Dudek, P. Häfner, S. Renaud, J. Schemmel, G. Cauwenberghs, J. Arthur, K. Hynna, F. Folowosele, S. SAA GHI, T. Serrano-Gotarredona, J. Wijekoon, Y. Wang, and K. Boahen. Neuromorphic silicon neuron circuits. *Frontiers in Neuroscience*, 5, 2011. ISSN 1662-453X. doi: 10.3389/fnins.2011.00073.
- N. Jain, S. Vaidyanath, A. Iyer, N. Natarajan, S. Parthasarathy, S. Rajamani, and R. Sharma. Jigsaw: Large language models meet program synthesis. *arXiv*, 2021.
- A. Q. Jiang, A. Sablayrolles, A. Mensch, C. Bamford, D. S. Chaplot, D. de las Casas, F. Bressand, G. Lengyel, G. Lample, L. Saulnier, L. R. Lavaud, M.-A. Lachaux, P. Stock, T. Le Scao, T. Lavril, T. Wang, T. Lacroix, and W. El Sayed. Mistral 7b. *arXiv preprint arXiv:2310.06825*, 2023.
- J. Johnson, M. Douze, and H. Jégou. Billion-scale similarity search with GPUs. *IEEE Transactions on Big Data*, 7(3):535–547, 2019.
- E. Jones and J. Steinhardt. Capturing failures of large language models via human cognitive biases. *arXiv preprint arXiv:2202.12299*, 2022.
- A. Karpathy. Software 2.0. Medium, 2017. URL <https://karpathy.medium.com/software-2-0-a64152b37c35>.
- N. Kassner, B. Krojer, and H. Schütze. Are Pretrained Language Models Symbolic Reasoners over Knowledge? In R. Fernández and T. Linzen (eds.), *Proceedings of the 24th Conference on Computational Natural Language Learning, CoNLL 2020, Online, November 19-20, 2020*, pp. 552–564. Association for Computational Linguistics, 2020. doi: 10.18653/v1/2020.conll-1.45.
- T. Katsch. Gateloop: Fully data-controlled linear recurrence for sequence modeling. *arXiv preprint arXiv:2311.01927*, 2023.
- D. Key, W.-D. Li, and K. Ellis. Toward trustworthy neural program synthesis. *arXiv preprint arXiv:2210.00848*, 2023.
- G. Kim, P. Baldi, and S. McAleer. Language models can solve computer tasks, 2023.
- L. Kirsch and J. Schmidhuber. Eliminating meta optimization through self-referential meta learning. *arXiv preprint arXiv:2212.14392*, 2022.

- A. Köpf, Y. Kilcher, D. von Rütte, S. Anagnostidis, Z.-R. Tam, K. Stevens, A. Barhoum, N. M. Duc, O. Stanley, R. Nagyfi, S. ES, S. Suri, D. Glushkov, A. Dantuluri, A. Maguire, C. Schuhmann, H. Nguyen, and A. Mattick. Openassistant conversations – democratizing large language model alignment, 2023.
- D. C. Kuncicky, S. I. Hruska, and R. C. Lacher. Hybrid systems: the equivalence of rule-based expert system and artificial neural network inference. *Int. J. Expert Syst.*, 4(3):281–297, jan 1991. ISSN 0894-9077.
- E. Kıcıman, R. Ness, A. Sharma, and C. Tan. Causal Reasoning and Large Language Models: Opening a New Frontier for Causality. *arXiv*, 2023.
- S. Lahlou, T. Deleu, P. Lemos, D. Zhang, A. Volokhova, A. Hernández-García, L. N. Ezzine, Y. Bengio, and N. Malkin. A theory of continuous generative flow networks. In *Proceedings of the International Conference on Machine Learning*, pp. 18269–18300. PMLR, 2023.
- J. E. Laird. Introduction to soar, 2022.
- J. E. Laird, A. Newell, and P. S. Rosenbloom. Soar: An architecture for general intelligence. *Artificial Intelligence*, 33(1):1–64, 1987. ISSN 0004-3702.
- L. C. Lamb, A. Garcez, M. Gori, M. Prates, P. Avelar, and M. Vardi. Graph neural networks meet neural-symbolic computing: A survey and perspective. In *AAAI Conference*, 2020.
- P. Langley, J. Laird, and S. Rogers. Cognitive architectures: Research issues and challenges. *Cognitive Systems Research*, 10:141–160, 2009. doi: 10.1016/j.cogsys.2006.07.004.
- Y. LeCun. A path towards autonomous machine intelligence, 2022. OpenReview Archive.
- H. Li, Y. Su, D. Cai, Y. Wang, and L. Liu. A survey on retrieval-augmented text generation. *arXiv preprint arXiv:2202.01110*, 2022a.
- Y. Li, D. Choi, J. Chung, N. Kushman, J. Schrittwieser, R. Leblond, T. Eccles, J. Keeling, F. Gimeno, A. Dal Lago, et al. Competition-level code generation with alphacode. *Science*, 378(6624):1092–1097, 2022b.
- Z. Li, Z. Yang, and M. Wang. Reinforcement learning with human feedback: Learning dynamic choices via pessimism. *arXiv preprint arxiv:2305.18438*, 2023.
- N. F. Liu, K. Lin, J. Hewitt, A. Paranjape, M. Bevilacqua, F. Petroni, and P. Liang. Lost in the middle: How language models use long contexts. *arXiv preprint arXiv:2307.03172*, 2023.
- S. Lu, D. Guo, S. Ren, J. Huang, A. Svyatkovskiy, A. Blanco, C. Clement, D. Drain, D. Jiang, D. Tang, et al. Codexglue: A machine learning benchmark dataset for code understanding and generation. *arXiv preprint arXiv:2102.04664*, 2021.
- M. Lutz. *Learning Python: Powerful Object-Oriented Programming*. Animal Guide. O’Reilly Media, 2013. ISBN 9781449355715.
- Q. Lyu, S. Havaldar, A. Stein, L. Zhang, D. Rao, E. Wong, M. Apidianaki, and C. Callison-Burch. Faithful chain-of-thought reasoning, 2023.
- J. MacGlashan, M. K. Ho, R. Loftin, B. Peng, G. Wang, D. L. Roberts, M. E. Taylor, and M. L. Littman. Interactive Learning from Policy-Dependent Human Feedback. In *Proceedings of the 34th International Conference on Machine Learning*, pp. 2285–2294. PMLR, July 2017.
- M. Macsweeney. Neural systems underlying british sign language and audio-visual english processing in native users. *Brain*, 125:1583–1593, 07 2002. doi: 10.1093/brain/awf153.
- A. Madaan, N. Tandon, P. Gupta, S. Hallinan, L. Gao, S. Wiegrefe, U. Alon, N. Dziri, S. Prabhume, Y. Yang, S. Welleck, B. P. Majumder, S. Gupta, A. Yazdanbakhsh, and P. Clark. Self-refine: Iterative refinement with self-feedback, 2023.
- K. Mahowald, A. A. Ivanova, I. A. Blank, N. Kanwisher, J. B. Tenenbaum, and E. Fedorenko. Dissociating language and thought in large language models: a cognitive perspective. *CoRR*, abs/2301.06627, 2023.

- R. Manhaeve, S. Dumancic, A. Kimmig, T. Demeester, and L. De Raedt. DeepProbLog: Neural Probabilistic Logic Programming. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018.
- J. Mao, C. Gan, P. Kohli, J. B. Tenenbaum, and J. Wu. The neuro-symbolic concept learner: Interpreting scenes, words, and sentences from natural supervision. In *7th International Conference on Learning Representations, ICLR 2019*, 2019.
- G. Marcus. The Next Decade in AI: Four Steps Towards Robust Artificial Intelligence. *arXiv preprint arXiv:2002.06177*, 2020.
- A. Martelli, A. Ravenscroft, and D. Ascher. *Python Cookbook*. O’Reilly Media, 2005. ISBN 9780596554743. URL <https://books.google.ro/books?id=Q0s6Vgb98CQC>.
- J. McCarthy. Lisp: A programming system for symbolic manipulations. In *Preprints of Papers Presented at the 14th National Meeting of the Association for Computing Machinery, ACM ’59*, pp. 1–4, New York, NY, USA, 1959. Association for Computing Machinery. ISBN 9781450373647. doi: 10.1145/612201.612243.
- J. McCarthy, M. L. Minsky, N. Rochester, and C. E. Shannon. A proposal for the dartmouth summer research project on artificial intelligence, august 31, 1955. *AI magazine*, 27(4):12–12, 2006.
- W. S. McCulloch and W. Pitts. A Logical Calculus of Ideas Immanent in Nervous Activity. *Bulletin of Mathematical Biophysics*, 5:115–133, 1943. doi: 10.1007/BF02478255.
- L. Menenti, S. M. E. Gierhan, K. Segaert, and P. Hagoort. Shared language: Overlap and segregation of the neuronal infrastructure for speaking and listening revealed by functional mri. *Psychological Science*, 22(9):1173–1182, 2011. doi: 10.1177/0956797611418347. PMID: 21841148.
- Microsoft. Bing is your AI-powered copilot for the web. Technical report, Microsoft, 2023. URL <https://bing.com/chat>.
- T. Mikolov, K. Chen, G. Corrado, and J. Dean. Efficient estimation of word representations in vector space, 2013a.
- Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space. *CoRR*, abs/1301.3781, 2013b. URL <http://dblp.uni-trier.de/db/journals/corr/corr1301.html#abs-1301-3781>.
- S. Milani, N. Topin, B. Houghton, W. H. Guss, S. P. Mohanty, K. Nakata, O. Vinyals, and N. S. Kuno. Retrospective analysis of the 2019 minerl competition on sample efficient reinforcement learning. In H. J. Escalante and R. Hadsell (eds.), *Proceedings of the NeurIPS 2019 Competition and Demonstration Track*, volume 123 of *Proceedings of Machine Learning Research*, pp. 203–214. PMLR, Dec 2020.
- L. De Moura and N. Bjørner. Z3: an efficient smt solver. In *Proceedings of the Theory and Practice of Software, 14th International Conference on Tools and Algorithms for the Construction and Analysis of Systems, TACAS’08/ETAPS’08*, pp. 337–340, Berlin, Heidelberg, 2008. Springer-Verlag. ISBN 3540787992.
- A. Newell. *Unified Theories of Cognition*. Harvard University Press, USA, 1990. ISBN 0674920996.
- A. Newell and H. Simon. The logic theory machine—a complex information processing system. *IRE Transactions on information theory*, 2(3):61–79, 1956.
- A. Newell and H. A. Simon. Human problem solving. *Prentice-Hall*, pp. 920, 1972.
- A. Newell and H. A. Simon. Computer science as empirical inquiry: symbols and search. *Commun. ACM*, 19(3): 113–126, mar 1976. ISSN 0001-0782. doi: 10.1145/360018.360022.
- A. Newell, J. C. Shaw, and H. A. Simon. Empirical explorations of the logic theory machine: a case study in heuristic. *IRE-AIEE-ACM ’57 (Western): Papers presented at the February 26-28, 1957, western joint computer conference: Techniques for reliability*, pp. 218–230, 1957. doi: 10.1145/1455567.1455605.
- H. Nori, Y. T. Lee, S. Zhang, D. Carignan, R. Edgar, N. Fusi, N. King, J. Larson, Y. Li, W. Liu, R. Luo, S. M. McKinney, R. O. Ness, H. Poon, T. Qin, N. Usuyama, C. White, and E. Horvitz. Can generalist foundation models outcompete special-purpose tuning? case study in medicine. *arXiv preprint arXiv:2311.16452*, 2023.

- A. v. d. Oord, S. Dieleman, H. Zen, K. Simonyan, O. Vinyals, A. Graves, N. Kalchbrenner, A. Senior, and K. Kavukcuoglu. Wavenet: A generative model for raw audio. *arXiv preprint arXiv:1609.03499*, 2016.
- OpenAI. Introducing ChatGPT. Technical report, OpenAI, November 2022. URL <https://openai.com/blog/chatgpt>.
- OpenAI. GPT-4 Technical Report. *arXiv*, 2023.
- L. Ouyang, J. Wu, X. Jiang, D. Almeida, C. L. Wainwright, P. Mishkin, C. Zhang, S. Agarwal, K. Slama, A. Ray, J. Schulman, J. Hilton, F. Kelton, L. E. Miller, M. Simens, A. Askell, P. Welinder, P. F. Christiano, J. Leike, and R. J. Lowe. Training language models to follow instructions with human feedback. *arXiv preprint arXiv:2203.02155*, 2022.
- F. Paischer, T. Adler, V. Patil, A. Bitto-Nemling, M. Holzleitner, S. Lehner, H. Eghbal-Zadeh, and S. Hochreiter. History compression via language models in reinforcement learning. In K. Chaudhuri, S. Jegelka, L. Song, C. Szepesvari, G. Niu, and S. Sabato (eds.), *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pp. 17156–17185. PMLR, July 2022.
- F. Paischer, T. Adler, M. Hofmarcher, and S. Hochreiter. Semantic helm: An interpretable memory for reinforcement learning. *CoRR*, abs/2306.09312, 2023. doi: 10.48550/arXiv.2306.09312.
- K. Papineni, S. Roukos, T. Ward, and W.-J. Zhu. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pp. 311–318, 2002.
- N. Park, D. Chae, J. Shim, S. Kim, E.-S. Kim, and J. Kim. Bridging the domain gap by clustering-based image-text graph matching. *arXiv preprint arXiv:2310.02692*, 2023.
- A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Köpf, E. Z. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala. Pytorch: An imperative style, high-performance deep learning library. *arXiv preprint arXiv:1912.01703*, abs/1912.01703, 2019.
- V. P. Patil, M. Hofmarcher, M.-C. Dinu, M. Dorfer, P. M. Blies, J. Brandstetter, J. A. Arjona-Medina, and S. Hochreiter. Align-RUDDER: Learning from few demonstrations by reward redistribution. *arXiv preprint arXiv:2009.14108*, 2020.
- F. Petroni, T. Rocktäschel, S. Riedel, P. S. H. Lewis, A. Bakhtin, Y. Wu, and A. H. Miller. Language Models as Knowledge Bases? In K. Inui, J. Jiang, V. Ng, and X. Wan (eds.), *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing, EMNLP-IJCNLP 2019, Hong Kong, China, November 3-7, 2019*, pp. 2463–2473. Association for Computational Linguistics, 2019. doi: 10.18653/v1/D19-1250.
- S. Pitis, M. R. Zhang, A. Wang, and J. Ba. Boosted prompt ensembles for large language models, 2023.
- M. Poli, S. Massaroli, E. Nguyen, D. Y. Fu, T. Dao, S. Baccus, Y. Bengio, S. Ermon, and C. Ré. Hyena hierarchy: Towards larger convolutional language models. *arXiv preprint arXiv:2302.10866*, 2023.
- B. Poole, A. Jain, J. T. Barron, and B. Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion. *arXiv preprint arXiv:2209.14988*, 2022.
- M. Qu and J. Tang. Probabilistic logic neural networks for reasoning. In *Proceedings of the 33rd International Conference on Neural Information Processing Systems*, 2019.
- R. Rafailov, A. Sharma, E. Mitchell, S. Ermon, C. D. Manning, and C. Finn. Direct preference optimization: Your language model is secretly a reward model. *arXiv preprint arXiv:2305.18290*, 2023.
- A. Ramesh, M. Pavlov, G. Goh, S. Gray, C. Voss, A. Radford, M. Chen, and I. Sutskever. Zero-shot text-to-image generation. In *International Conference on Machine Learning*, pp. 8821–8831. PMLR, 2021.
- H. Ramsauer, B. Schäfl, J. Lehner, P. Seidl, M. Widrich, T. Adler, L. Gruber, M. Holzleitner, M. Pavlović, G. K. Sandve, et al. Hopfield networks is all you need. *arXiv preprint arXiv:2008.02217*, 2020.
- M. Regev, C. J. Honey, E. Simony, and U. Hasson. Selective and invariant neural responses to spoken and written narratives. *Journal of Neuroscience*, 33(40):15978–15988, 2013. ISSN 0270-6474. doi: 10.1523/JNEUROSCI.1580-13.2013.

- ReplikaAI. Pushing the Boundaries of AI to Talk to the Dead. Technical report, ReplikaAI, 2016. URL <https://www.bloomberg.com/news/articles/2016-10-20/pushing-the-boundaries-of-ai-to-talk-to-the-dead>.
- B. Romera-Paredes, M. Barekatin, A. Novikov, et al. Mathematical discoveries from program search with large language models. *Nature*, 2023. doi: 10.1038/s41586-023-06924-6.
- F. Rosenblatt. The perceptron: A probabilistic model for information storage and organization in the brain. *Psychological Review*, 65(6):386–408, 1958. doi: 10.1037/h0042519.
- L. Ruis, A. Khan, S. Biderman, S. Hooker, T. Rocktäschel, and E. Grefenstette. Large language models are not zero-shot communicators. *CoRR*, abs/2210.14986, 2022. doi: 10.48550/arXiv.2210.14986.
- D. Rumelhart, G. Hinton, and R. Williams. Learning representations by back-propagating errors. *Nature*, 323:533–536, 1986. doi: 10.1038/323533a0.
- C. Saharia, W. Chan, S. Saxena, L. Li, J. Whang, E. Denton, S. K. S. Ghasemipour, B. K. Ayan, S. S. Mahdavi, R. G. Lopes, et al. Photorealistic text-to-image diffusion models with deep language understanding. *arXiv preprint arXiv:2205.11487*, 2022.
- A. Santoro, A. Lampinen, K. Mathewson, T. Lillicrap, and D. Raposo. Symbolic behaviour in artificial intelligence. *arXiv preprint arXiv:2102.03406*, 2022.
- T. Schick, J. Dwivedi-Yu, R. Dessi, R. Raileanu, M. Lomeli, L. Zettlemoyer, N. Cancedda, and T. Scialom. Toolformer: Language models can teach themselves to use tools, 2023.
- J. Schmidhuber. Gödel machines: Fully self-referential optimal universal self-improvers. *Cognitive Technologies*, 8: 199–226, 01 2007. doi: 10.1007/978-3-540-68677-4\_7.
- J. Schmidhuber. Driven by compression progress: A simple principle explains essential aspects of subjective beauty, novelty, surprise, interestingness, attention, curiosity, creativity, art, science, music, jokes. *arXiv preprint arXiv:0812.4360*, 2009.
- J. Schrittwieser, I. Antonoglou, T. Hubert, K. Simonyan, L. Sifre, S. Schmitt, A. Guez, E. Lockhart, D. Hassabis, T. Graepel, T. Lillicrap, and D. Silver. Mastering atari, go, chess and shogi by planning with a learned model. *Nature*, 588(7839):604–609, 2020. doi: 10.1038/s41586-020-03051-4.
- K. Schweighofer, A. Radler, M.-C. Dinu, M. Hofmarcher, V. P. Patil, A. Bitto-Nemling, H. Eghbal-zadeh, and S. Hochreiter. A dataset perspective on offline reinforcement learning. In *Conference on Lifelong Learning Agents*, pp. 470–517. PMLR, 2022.
- T. Scott, J. Gallée, and E. Fedorenko. A new fun and robust version of an fMRI localizer for the frontotemporal language system. *Cognitive Neuroscience*, 8:1–10, 07 2016. doi: 10.1080/17588928.2016.1201466.
- N. Shinn, B. Labash, and A. Gopinath. Reflexion: an autonomous agent with dynamic memory and self-reflection, 2023.
- D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, G. van den Driessche, J. Schrittwieser, I. Antonoglou, V. Panneershelvam, M. Lanctot, S. Dieleman, D. Grewe, J. Nham, N. Kalchbrenner, I. Sutskever, T. P. Lillicrap, M. Leach, K. Kavukcuoglu, T. Graepel, and D. Hassabis. Mastering the game of Go with deep neural networks and tree search. *Nature*, 529(7587):484–489, 2016. doi: 10.1038/nature16961.
- D. Silver, T. Hubert, J. Schrittwieser, I. Antonoglou, M. Lai, A. Guez, M. Lanctot, L. Sifre, D. Kumaran, T. Graepel, T. Lillicrap, K. Simonyan, and D. Hassabis. Mastering chess and shogi by self-play with a general reinforcement learning algorithm. *arXiv preprint arXiv:1712.01815*, 2017a.
- D. Silver, J. Schrittwieser, K. Simonyan, et al. Mastering the game of go without human knowledge. *Nature*, 550: 354–359, 2017b. doi: 10.1038/nature24270.
- U. Singer, A. Polyak, T. Hayes, X. Yin, J. An, S. Zhang, Q. Hu, H. Yang, O. Ashual, O. Gafni, et al. Make-a-video: Text-to-video generation without text-video data. *arXiv preprint arXiv:2209.14792*, 2022.
- K. Singhal, S. Azizi, T. Tu, S. S. Mahdavi, J. Wei, H. W. Chung, N. Scales, et al. Large language models encode clinical knowledge. *Nature*, 620(7972):172–180, 2023.

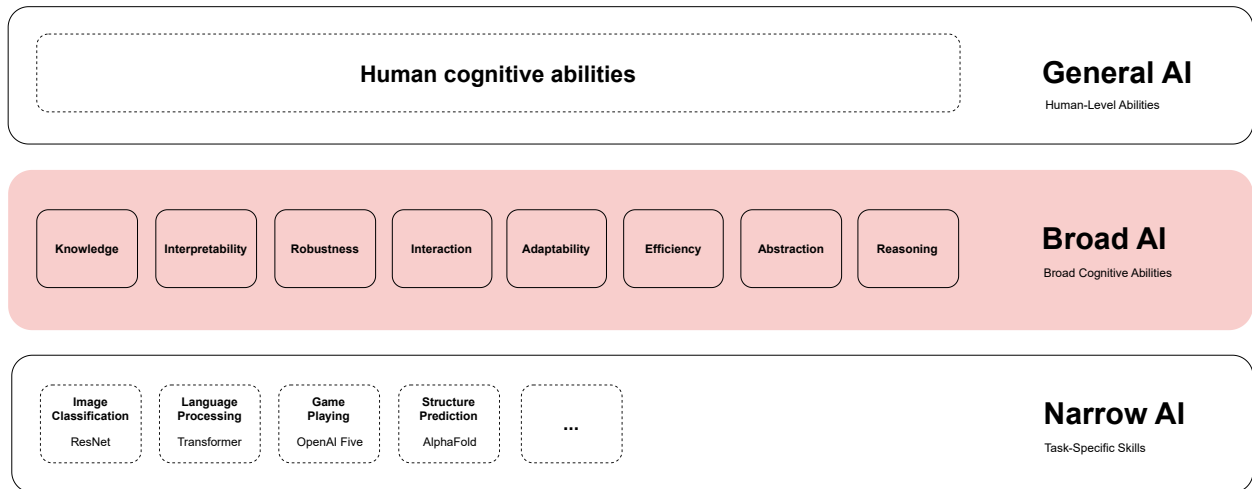


- K. Song, X. Tan, T. Qin, J. Lu, and T.-Y. Liu. MpNet: Masked and permuted pre-training for language understanding. In *Proceedings of the 34th International Conference on Neural Information Processing Systems, NIPS'20*, pp. 1414, Red Hook, NY, USA, 2020. Curran Associates Inc. ISBN 9781713829546.
- Spotify. Approximate Nearest Neighbors Oh Yeah. Technical report, Spotify, 2017.
- T. R. Sumers, S. Yao, K. Narasimhan, and T. L. Griffiths. Cognitive architectures for language agents. *arXiv preprint arXiv:2309.02427*, 2023.
- M. Summerfield. *Programming in Python 3: A Complete Introduction to the Python Language*. Developer's library. Addison-Wesley, 2010. ISBN 9780321680563.
- J. J. Sun, M. Tjandrasuwita, A. Sehgal, A. Solar-Lezama, S. Chaudhuri, Y. Yue, and O. Costilla-Reyes. Neurosymbolic programming for science. *arXiv preprint arXiv:2210.05050*, 2022.
- R. S. Sutton. *Temporal Credit Assignment in Reinforcement Learning*. PhD thesis, University of Massachusetts, Dept. of Comp. and Inf. Sci., 1984.
- P. Swazinna, S. Udluft, D. Hein, and T. Runkler. Comparing model-free and model-based algorithms for offline reinforcement learning. *arXiv preprint arXiv:2201.05433*, 2022.
- Z. Szabó, B. K. Sriperumbudur, B. Póczos, and A. Gretton. Learning theory for distribution regression. *J. Mach. Learn. Res.*, 17(1):5272–5311, Jan 2016. ISSN 1532-4435.
- R. Taori, I. Gulrajani, T. Zhang, Y. Dubois, X. Li, C. Guestrin, P. Liang, and T. B. Hashimoto. Stanford alpaca: An instruction-following llama model. [https://github.com/tatsu-lab/stanford\\_alpaca](https://github.com/tatsu-lab/stanford_alpaca), 2023.
- H. Touvron, T. Lavril, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, B. Rozière, N. Goyal, E. Hambro, F. Azhar, A. Rodriguez, A. Joulin, E. Grave, and G. Lample. Llama: Open and efficient foundation language models, 2023.
- L. Tunstall, E. Beeching, N. Lambert, N. Rajani, K. Rasul, Y. Belkada, S. Huang, L. von Werra, C. Fourier, N. Habib, N. Sarrazin, O. Sanseviero, A. M. Rush, and T. Wolf. Zephyr: Direct distillation of lm alignment, 2023.
- A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin. Attention is all you need. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.
- Ramakrishna Vedantam, C. Lawrence Zitnick, and Devi Parikh. Cider: Consensus-based image description evaluation. *CoRR*, abs/1411.5726, 2014. URL <http://arxiv.org/abs/1411.5726>.
- P. Veličković and C. Blundell. Neural algorithmic reasoning. *Patterns*, 2(7):100273, 2021. ISSN 2666-3899. doi: <https://doi.org/10.1016/j.patter.2021.100273>.
- B. Wang, Z. Wang, X. Wang, Y. Cao, R. A. Saurous, and Y. Kim. Grammar prompting for domain-specific language generation with large language models. *arXiv preprint arXiv:2305.19234*, 2023a.
- J. Wang, X. Yi, R. Guo, H. Jin, P. Xu, S. Li, X. Wang, X. Guo, C. Li, X. Xu, et al. Milvus: A purpose-built vector data management system. In *Proceedings of the 2021 International Conference on Management of Data*, pp. 2614–2627, 2021a.
- X. Wang, J. Wei, D. Schuurmans, Q. Le, E. Chi, S. Narang, A. Chowdhery, and D. Zhou. Self-consistency improves chain of thought reasoning in language models. *arXiv preprint arXiv:2203.11171*, 2023b.
- Y. Wang, R.J. Skerry-Ryan, D. Stanton, Y. Wu, R. J. Weiss, N. Jaitly, Z. Yang, Y. Xiao, Z. Chen, S. Bengio, et al. Tacotron: Towards end-to-end speech synthesis. *arXiv preprint arXiv:1703.10135*, 2017.
- Y. Wang, W. Wang, S. Joty, and S. C. H. Hoi. Codet5: Identifier-aware unified pre-trained encoder-decoder models for code understanding and generation. *arXiv preprint arXiv:2109.00859*, 2021b.
- Y. Wang, Y. Kordi, S. Mishra, A. Liu, N. A. Smith, D. Khashabi, and H. Hajishirzi. Self-instruct: Aligning language model with self generated instructions. *arXiv preprint arXiv:2212.10560*, 2022.
- J. Wei, Y. Tay, R. Bommasani, C. Raffel, B. Zoph, S. Borgeaud, D. Yogatama, M. Bosma, D. Zhou, D. Metzler, E. H. Chi, T. Hashimoto, O. Vinyals, P. Liang, J. Dean, and W. Fedus. Emergent abilities of large language models. *arXiv preprint arXiv:2206.07682*, 2022a.

- J. Wei, X. Wang, D. Schuurmans, M. Bosma, B. Ichter, F. Xia, E. H. Chi, Q. V. Le, and D. Zhou. Chain of thought prompting elicits reasoning in large language models. In A. H. Oh, A. Agarwal, D. Belgrave, and K. Cho (eds.), *Advances in Neural Information Processing Systems*, 2022b.
- Y. Weng, M. Zhu, F. Xia, B. Li, S. He, S. Liu, B. Sun, K. Liu, and J. Zhao. Large language models are better reasoners with self-verification. *arXiv preprint arXiv:2212.09561*, 2023.
- A. N. Whitehead and B. Russell. *Principia Mathematica*. Cambridge University Press, 1925–1927.
- M. Widrich, M. Hofmarcher, V. P. Patil, A. Bitto-Nemling, and S. Hochreiter. Modern Hopfield Networks for Return Decomposition for Delayed Rewards. In *Deep RL Workshop NeurIPS 2021*, 2021.
- Writesonic. ChatGPT Alternative Built With Superpowers - ChatSonic. Technical report, Chatsonic, 2022. URL <https://writesonic.com/chat>.
- C. Xu, D. Guo, N. Duan, and J. McAuley. Baize: An open-source chat model with parameter-efficient tuning on self-chat data, 2023.
- Z. Xu, H. van Hasselt, and D. Silver. Meta-gradient reinforcement learning. *ArXiv*, 2018.
- L. Yang, S. Zhang, Z. Yu, G. Bao, Y. Wang, J. Wang, R. Xu, W. Ye, X. Xie, W. Chen, and Y. Zhang. Supervised Knowledge Makes Large Language Models Better In-context Learners. *arXiv preprint arXiv:2312.15918*, 2023.
- S. Yao, D. Yu, J. Zhao, I. Shafran, T. L. Griffiths, Y. Cao, and K. Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. *arXiv preprint arXiv:2305.10601*, 2023a.
- S. Yao, J. Zhao, D. Yu, N. Du, I. Shafran, K. Narasimhan, and Y. Cao. React: Synergizing reasoning and acting in language models. *arXiv preprint arXiv:2210.03629*, 2023b.
- J. Ye, Z. Wu, J. Feng, T. Yu, and L. Kong. Compositional exemplars for in-context learning. *arXiv preprint arXiv:2302.05698*, 2023.
- K. You, X. Wang, M. Long, and M. Jordan. Towards accurate model selection in deep unsupervised domain adaptation. In K. Chaudhuri and R. Salakhutdinov (eds.), *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pp. 7124–7133. PMLR, Jun 9–15 2019.
- YouWrite. The AI Search Engine You Control. Technical report, You.com, 2022. URL <https://you.com>.
- D. Yu, B. Yang, D. Liu, H. Wang, and S. Pan. A survey on neural-symbolic learning systems. *Neural Networks*, 166: 105–126, 2023. ISSN 0893-6080.
- W. Yuan, R. Y. Pang, K. Cho, S. Sukhbaatar, J. Xu, and J. Weston. Self-rewarding language models. *arXiv preprint arXiv:2401.10020*, 2024.
- J. Zhang, B. Chen, L. Zhang, X. Ke, and H. Ding. Neural, symbolic and neural-symbolic reasoning on knowledge graphs. *AIOpen*, pp. 14–35, 2021.
- M. Zhuge, W. Wang, L. Kirsch, F. Faccio, D. Khizbullin, and J. Schmidhuber. Language agents as optimizable graphs, 2024.

## A BROAD AI AND NEURO-SYMBOLIC SYSTEMS

Our work focuses on broad AI (Hochreiter, 2022) (see Figure 7) through the integration of symbolic and sub-symbolic AI methodologies. Broad AI extends beyond restricted focus on single-task performance of narrow AI. In broad AI, systems are engineered to handle a range of tasks with a high degree of autonomy, utilizing sensory input, accumulated experiences, and previously developed skills.



**Figure 7:** Hierarchical model of “cognitive” abilities of AI systems (Chollet, 2019; Hochreiter, 2022). The figure contrasts the emergent paradigm of *Broad AI* with current *Narrow AI* systems, showcasing Broad AI’s wider range of capabilities, such as knowledge transfer, interaction, adaptability, robustness, abstraction, advanced reasoning, and efficiency. Broad AI aims to mimic human cognitive adaptability and robustness through advanced methodologies like few-shot learning, self-supervised contrastive learning, and context-sensitive sensory processing. Notably, Broad AI applies principles such as conceptual short-term memory and modern Hopfield networks (Ramsauer et al., 2020) to better integrate context and memory, thus avoiding pitfalls like explaining away and short-cut learning. We acknowledge the potential of NeSy systems as a significant step towards AI systems capable of performing any cognitive task with human-like proficiency.

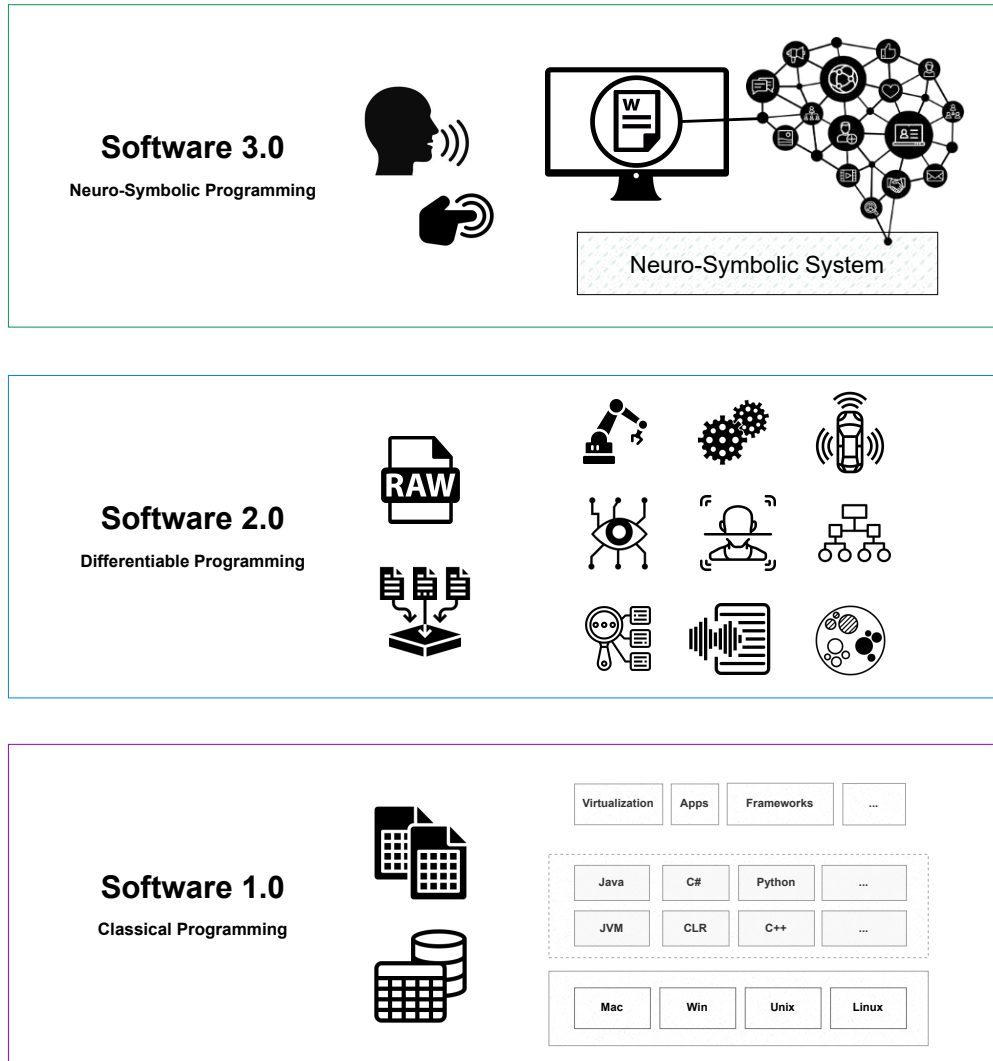
NeSy methods form the basis for developing new cognitive architectures (Newell & Simon, 1956; Newell et al., 1957; Newell & Simon, 1972; Newell, 1990; Langley et al., 2009; Laird, 2022; Dawid & LeCun, 2023; Summers et al., 2023; LeCun, 2022; Assran et al., 2023). This hybridization produces computational graphs capable of context-aware learning and reasoning, allowing AI to execute complex tasks with human-like flexibility.

Borrowing nomenclature from Karpathy (2017); Dilhara et al. (2021), we refer to the next generation of software as Software 3.0, which consists of applications that are neither pre-determined at design time, nor learned through statistical inference, but triggered by an interaction which stimulates the realization of a computational graph analogous to *neuromorphic circuits* (Indiveri et al., 2011), however, purely established at inference time in the “*thought*” process of a NeSy system.

To enable such systems, we require a more native integration (see illustration in Figure 9) of probabilistic programming paradigms into our contemporary programming stack, and make their utilization a commodity for practitioners and researchers alike.

### A.1 BROADER IMPACT

With LLMs becoming more and more accessible, progress recently made possible by the vast open source contributions from Köpf et al. (2023); Touvron et al. (2023); Taori et al. (2023); Xu et al. (2023); Geng et al. (2023); Biderman et al. (2023), embedded accelerators for LLMs — or more generally NeSy engines — will be ubiquitous in future computation platforms, such as wearables, smartphones, tablets, consoles, or notebooks. Although current execution cycles are slow and error-prone, we expect to see further performance gains through improved operating system level optimizations, dedicated GPU-centric hardware refinement, and improved software interoperability. We believe that modern programming paradigms should natively support probabilistic concepts and provide a boilerplate-free set of features for constructing and evaluating generative computational graphs. This includes but is not limited to compo-

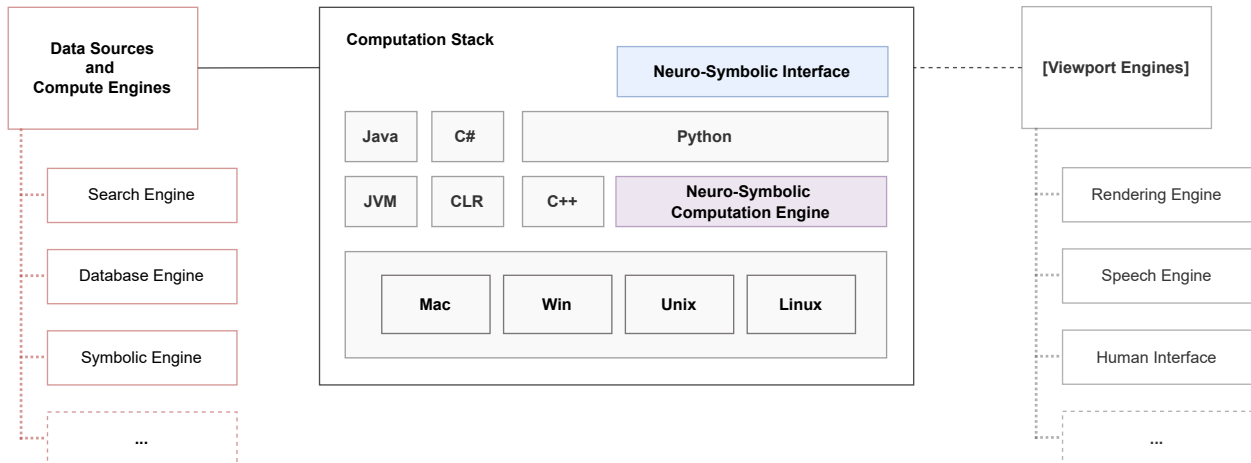


**Figure 8:** Evolution of software paradigms: From Software 1.0’s rigid specification in classical programming to Software 2.0’s data-driven and objective function-focused differentiable programming, leading to Software 3.0’s NeSy systems that emphasize human-centric, interaction-based programming with computational graphs. This progression represents a shift from explicit task-specific programming to abstract, adaptive systems that cater to dynamic user preferences.

sitional, parallelizable, and simulation-based executions with polymorphic operations and self-referential structures. Current programming languages often have disjointed or makeshift solutions for these concepts in the context of generative processes. We believe integral probabilistic support for these concepts into modern software and hardware will unlock new programming paradigms that can fully take advantage of generative architectures. We hope the community will consider these ideas as essential components of contemporary computing.

We also expect to see significant progress by processing central language concepts through system-on-a-chip (SoC) solutions of pre-trained models, with linear probing layers for hot-swappable weight exchange of task-specific projections and executions. A wide range of functionalities can be then offloaded to probabilistic programming languages to operate on dedicated symbols and streamline the vector-valued mappings between the concept space and underlying problem space, avoiding defining boilerplate code to load and unload network weights.

Furthermore, we believe that many gains in representational stability and consistency may be obtained through multi-modal data training and improved alignment based on operator learning oriented functionalities and workflow-related scoring functionalities, analogous to our introduced quality measure. Gains in representational stability also benefit



**Figure 9:** The illustration demonstrates the integration of Neuro-Symbolic computation within the contemporary programming stack. Probabilistic programming paradigms are embedded natively alongside traditional languages and environments, facilitated by interfaces to various data sources, compute engines, and human interaction tools, streamlining their adoption in practical and research applications.

self-instruction and self-referential sub-process evaluations, which enable the dynamic creation and evaluation of complex hierarchical computational graphs. This will enable online learning models to perform, in real-time, skill acquisition of complex concepts with only one or few examples at inference time. We believe this will enable the creation of autonomously self-evolving cognitive architectures (Langley et al., 2009; Dawid & LeCun, 2023; Sumers et al., 2023). We therefore see an inherent connection to generative design as an analogy for creating coherent and stable “*thought*” computational graphs, and believe this paves the path toward broad AI systems (see Section A) and is a requirement for developing artificial general intelligent agents.

Finally, we also wish to express our concern about recent economic trends in the deep-tech industry, where we observe AI-related concentration of data and resources, coupled with a tendency towards closed-source practices. We strongly advocate for increased transparency and exchange of ideas to ensure a diverse and collective growth in our socio-economic landscape. Therefore, we push towards a democratic and open-source initiative.

## B CONNECTION BETWEEN FRÉCHET DISTANCE AND MAXIMUM MEAN DISCREPANCY

Let us consider a Gaussian kernel defined by the expression

$$K(x, y) = \exp\left(-\frac{\|x - y\|^2}{2\sigma^2}\right), \quad (5)$$

where  $\sigma$  is the bandwidth parameter of the kernel and  $\|\cdot\|$  denotes the Euclidean norm. Utilizing  $K$ , we can now construct a measure of distance between distributions, by embedding them into the Reproducing Kernel Hilbert Space (RKHS) induced by  $K$ , using kernel mean embeddings. The resulting distance is called the Maximum Mean Discrepancy (MMD).

More precisely, the MMD between two probability distributions  $P$  and  $Q$  is encoded in the RKHS through mean embeddings, which can be expressed as

$$\text{MMD}^2(P, Q) = \|\mathbb{E}_{x \sim P}[\phi(x)] - \mathbb{E}_{y \sim Q}[\phi(y)]\|_{\text{RKHS}}^2, \quad (6)$$

where  $\phi(\cdot)$  represents the feature mapping to the RKHS corresponding to the kernel  $K$ .

On the other hand, for multivariate Gaussian distributions, we can use the Fréchet distance as a measure of similarity, which is nothing but the associated Wasserstein-2 distance, for which an explicit formula is available in the Gaussian case. The resulting expression is as follows (Dowson & Landau, 1982):

$$d^2(X_1, X_2) = \|\mu_1 - \mu_2\|_2^2 + \text{Tr}\left(C_1 + C_2 - 2(C_1 C_2)^{\frac{1}{2}}\right), \quad (7)$$

where  $X_1 \sim \mathcal{N}(\mu_1, C_1)$  and  $X_2 \sim \mathcal{N}(\mu_2, C_2)$ , and  $\text{Tr}(\cdot)$  indicates the trace of a matrix.

To establish an approximation of the Fréchet distance using the Gaussian kernel, we take  $C_1 = \sigma^2 I$  and  $C_2 = \sigma^2 I$  as identity covariance matrices scaled by  $\sigma^2$ . This assumption allows us to focus solely on the disparities in mean vectors:

$$d^2(X_1, X_2) \approx \|\mu_1 - \mu_2\|_2^2, \quad (8)$$

setting aside the effect of different covariance structures.

Given these conditions, we attempt to argue that the Fréchet distance behaves similarly as MMD:

$$d^2(X_1, X_2) \approx \|\mu_1 - \mu_2\|_2^2 \approx \text{MMD}^2(P, Q), \quad (9)$$

Heuristically, at least for small  $\|\mu_1 - \mu_2\|$ , also the associated kernel evaluations  $K(X_1, X_2)$  tend to be small (see also [Hochreiter & Schmidhuber \(1997\)](#)), which leads to a small MMD, if we ignore the terms associated to  $K(X_1, X_1)$ ,  $K(X_2, X_2)$  (which cancel out due to same covariance structure).

In the next section, we want to further elaborate on the MMD and a possible score, that can be derived from it.

### B.1 EXTENDED SIMPLIFICATION OF THE MMD CALCULATION

To understand the simplification of the MMD, we are formally expressing the MMD in terms of kernel sums over pairs of samples within and across two samples  $X$  and  $Y$ :

$$\text{MMD}^2(X, Y) = \frac{1}{m(m-1)} \sum_i \sum_{j \neq i} k(x_i, x_j) - \frac{2}{mn} \sum_{i=1}^m \sum_{j=1}^n k(x_i, y_j) + \frac{1}{n(n-1)} \sum_i \sum_{j \neq i} k(y_i, y_j), \quad (10)$$

where  $m$  and  $n$  are the sizes of samples  $X$  and  $Y$ , respectively.

Empirical observations have led to the conclusion that the within-sample terms  $\sum_i \sum_{j \neq i} k(x_i, x_j)$  and  $\sum_i \sum_{j \neq i} k(y_i, y_j)$  cancel out the cross terms  $\sum_{i=1}^m \sum_{j=1}^n k(x_i, y_j)$  under certain conditions. This can be due to the following:

- In high-dimensional embedding spaces, distributions of embedding vectors are often closely related and normally distributed.
- If the samples  $X$  and  $Y$  are drawn from distributions  $P$  and  $Q$  where their mean embeddings are nearly orthogonal in the RKHS, it is the dissimilarity across samples, rather than that within, that is most relevant.

Therefore, under these specific conditions, it becomes justifiable to focus on the cross-terms, yielding the following proposal for a distance measure:

$$\widetilde{\text{MMD}}^2(X, Y) \approx \frac{2}{mn} \sum_{i=1}^m \sum_{j=1}^n k(x_i, y_j). \quad (11)$$

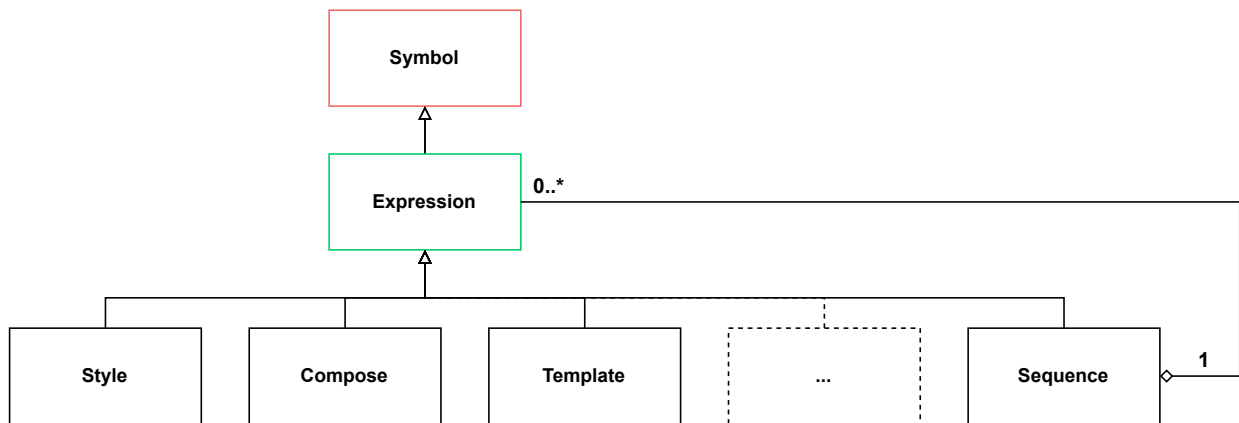
## C STRUCTURE

**Primitives** In the SymbolicAI framework, at the core lies the concept of Primitives and the dynamic type creation of `Symbol` objects, which are central to inherit types of behaviors. Primitives are pre-defined operations that act on `Symbol` objects, encapsulating basic operations, such as arithmetic, logic, or casting operations, to name a few. These operations are crucial to the framework’s versatility and form the foundation for more complex interactions within computational graphs. Essentially, they can be viewed as contextualized functions that accept a `Symbol` object, send it to the NeSy engine for evaluation, and return one or more new objects, primarily new symbols. One of the key features of operations is their polymorphism, which allows for them to be applied to various data types, such as strings, integers, floats, lists, and more, with different behaviors depending on the specific object instance. To execute operations, we utilize the `Symbol` object’s `value` attribute containing the original data type. This will be then sent as a string representation to the engines to execute the needed operations. Consequently, all values are cast to a string representation. Remember, this was our implicit assumption (see Section 4). For custom objects, it is essential to define a suitable `__str__` method to cast the object to a string representation while preserving the object’s semantics.

**Symbol Objects Creation and Dynamic Typing** A `Symbol` object is a versatile entity that can encapsulate a variety of data types and behaviors. The creation of `Symbol` objects is facilitated through a metaclass, which enables the dynamic typing of these objects to inherit behaviors from a collection of primitives. This dynamic typing system is important for extending the functionality of `Symbol` objects beyond simple data containers; they contain specific behaviors appropriate for the operations they will perform. For instance, a `Symbol` object may possess the behaviors of arithmetic computations, string manipulations, or even logical comparisons, depending on the defined primitives.

**Type Inheritance and Expression Creation** Type inheritance in SymbolicAI is leveraged to create new expressions, which are specialized forms of `Symbol` objects designed to represent parts of a computational graph. Expressions extend the capabilities of `Symbol` objects by providing a structured way to define complex functionalities that can later be evaluated to produce new `Symbol` objects or modify existing ones. In SymbolicAI, an `Expression` object inherits the properties of `Symbol` objects while also being able to define its own unique behavior through a `forward` method, which is analogous to a computational graph node’s evaluation function. Figure 10 gives an overview of an exemplary inheritance branch. Each `Expression` must feature a `forward` method, which must be overwritten to define its behavior. The inherited `__call__` method invokes the `forward` method, evaluating the expression and returning the result. This design pattern facilitates lazy evaluation of expressions, allowing for complex composition of expressions.

Inherited from the `Symbol` class, the `_sym_return_type` and `static_context` properties establish the context in which the current `Expression` operates. The `static_context` impacts all operations of the current `Expression` subclass, while the `_sym_return_type` guarantees the acquisition of the desired return object type post-evaluation. Typically, this returns the current type but can be configured to return a different type. A more in-depth examination of both notions will be provided in the following section.



**Figure 10:** Class diagram showing the inheritance and composition relationships among `Symbol`, `Expression`, and other inherited classes. `Symbol` serves as a base class for `Expression` where all the other types are derived from. Other types may contain or associate with zero or more `Symbol` types. For example, we illustrate how the `Sequence` derives from `Expression` and the multiplicity `'0..*'` indicates that a `Sequence` can contain any number of `Expression` instances.

**Utilizing Decorators for Operation Definition** Decorators serve as a bridge between the declarative nature of symbolic operations and the imperative execution model of programming languages. By augmenting function definitions with decorators, the framework can dynamically assign operations to `Symbol` or `Expression` objects, which are then interpreted by the underlying NeSy engine or traditional solvers.

For example, the `@core.logic` decorator can be used to augment a `Symbol` object with the capability to perform logical `and`, `or`, or `not` operations contextually. Similarly, the `@core.combine` decorator allows the framework to define the semantics of combining or adding two symbolic values, regardless of their underlying data representations.

```

1 # Example of using decorators to define logical operations
2 @core.logic(operator='and')
3 def _some_logic(self, other):
4     # implementation logic here
5     pass

```

**Aspect-Oriented Programming** The aspect-oriented programming paradigm offers a functional approach for extending or modifying the behavior of functions or methods without altering their code directly. This adheres to the principles of modularity and separation of concerns, as it allows for the isolation of specific functionalities while maintaining the original function’s core purpose. By wrapping the original function, decorators provide an efficient and reusable way of adding or modifying behaviors. For instance, SymbolicAI integrates the zero- and few-shot learning with default fallback functionalities of pre-existing code.

Decorators brings several advantages (Beazley, 2009; Martelli et al., 2005; Summerfield, 2010; Lutz, 2013):

- **Reusability:** Decorators promote code modularity, enhancing code reusability and contributing to software maintainability. This advantage is particularly salient when managing a variety of operations, reducing redundancy and simplifying the integration of new functionalities.
- **Composition:** Decorators support function composition, allowing developers to construct complex functionalities from pre-existing code blocks without the need to expand the codebase or rely on complex inheritance hierarchies.
- **Adaptability:** Through decorators we can easily modify or extend the behavior of operations without changing their core implementation. This flexibility facilitates the generation of adaptive workflows and reliable fallback mechanisms when experimental implementations do not fulfill required constraints.

**Symbol Class and Computational Graph Elements** A computational graph in the SymbolicAI framework is an assembly of interconnected `Symbol` objects, each encapsulating a unit of data and the operations that can be performed on it. The exchange between these symbols forms a highly modular and interpretable system, capable of representing complex workflows.

The `Symbol` class is an abstraction representing data and context. It holds not only the value itself, but metadata that guides its transformation and interpretation. Through inheritance and compositionality, the `Symbol` can be extended into more complex expressions, and becoming nodes in a computational graph. Each `Symbol` instance can optionally contain a reference to its parent and children, naturally forming a directed graph structure where the nodes are symbols and edges represent relationships between a symbol and its derivative computations.

The `Linker` class, is a metadata subclass, and tracks relationships and results, effectively annotating the graph with execution details. It keeps records of nodes’ keys, allowing quick retrieval of related computational outcomes within the graph, and aids in tasks such as error tracing and debugging.

A central concept in this structure is the notion of `root`, which points to the origin of the computational sequence. Accessing the root allows backtracking through the graph, making it possible to aggregate results and inspect the flow of computation that led to the current node.

The computational graph’s structure is further enriched by properties like `nodes`, `edges`, and `graph` itself, which collectively enable the comprehensive query of the computation’s topology. These properties are used internally to enable graph visualizations, which are useful for debugging and analysis.

**Expression of a Computational Graph** In practice, consider the `Expression` class, which extends the functionality of the `Symbol` class. When composing a `Sequence` of `Expression` objects, we are effectively composing operations in a predetermined order.

For instance, an expression like:

```
1 Sequence (
2   Clean(),
3   Translate(),
4   Outline(),
5   Compose(),
6 )
```

represents a procedure that first cleans data, then translates it, outlines the essential information, and composes it into a finalized form. When this sequence is executed, the operations unfold in the exact order specified, with each step receiving the output of its predecessor as input and passing its result to the successor.



## D INSTALLATION

The installation of the SymbolicAI framework is straightforward and can be done through the Python package manager pip. To install SymbolicAI, open a terminal and execute the following command in your current python environment:

```
1 pip install symbolicalai
```

This command will install the latest version of SymbolicAI along with its core dependencies, enabling the integration of the framework into Python applications. If you intend to utilize the framework with local engines<sup>5</sup>, or with engines powered by external APIs such as OpenAI’s API, additional installation steps are required.

### D.1 ENGINE CONFIGURATION

Before the first run, it is necessary to configure the required modules and optionally set necessary API keys to activate the respective engines. This can be done in multiple ways, but we recommend doing it through the configuration wizard by running this command in the terminal:

```
1 symwzd
```

This step is essential to register the engines internally for subsequent runs.

For instance, SymbolicAI includes OpenAI’s GPT models as NeSy engine. To only set or change OpenAI API keys, the following command is used before starting a SymbolicAI instance:

```
1 # Linux / MacOS
2 export OPENAI_API_KEY="<OPENAI_API_KEY>"
```

After setting up the API keys, the SymbolicAI library is imported in Python as following:

```
1 import symai
```

For more low-level changes, we store everything under the `$HOME/.symai` folder, such as the `symai.config.json`, which stores every key, both registered and not registered.

### D.2 OPTIONAL INSTALLATIONS

The SymbolicAI framework is designed to leverage multiple engines for a variety of operations. To fully utilize these capabilities, you may install additional dependencies or set up the optional API keys for specific engines like WolframAlpha, SerpApi, and others. In Figure 11 we conceptually outline the connection between the utilization of an LLM and its interact with other tools and solvers. Instructions and operations can be initiated by any user, pre-scripted knowledge base or learned meta agent.

For instructions on additional installations, including the support of optional engines, refer to the documentation provided with the framework. This documentation will give detailed steps on installing optional dependencies and configuring additional API keys.

## E IMPLEMENTATION DETAILS

Let us now define some `Symbol` objects and perform some basic manipulations.

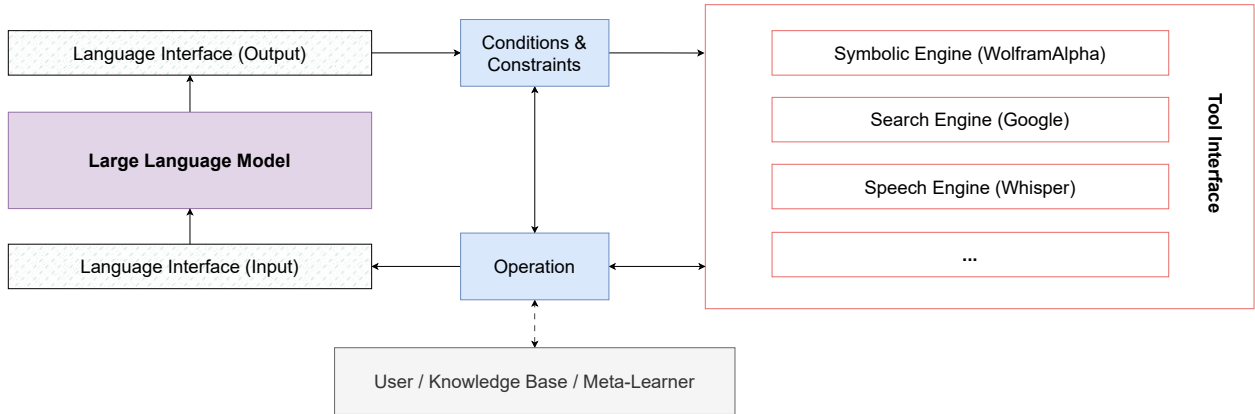
### E.1 FUZZY COMPARISON

For instance, let’s consider fuzzy<sup>6</sup> comparisons. Within SymbolicAI, it enables more adaptable and context-aware evaluations, accommodating the inherent uncertainties and variances often encountered in real-world data.

```
1 import numpy
2
3 s = symai.Symbol('3.1415...')
4 s == numpy.pi
```

<sup>5</sup> The local engines are experimental and are run on your local machine. For more details, refer to the “Local Neuro-Symbolic Engine” section in the documentation.

<sup>6</sup> Not related to fuzzy logic, which is a topic under active consideration.



**Figure 11:** The SymbolicAI framework integrates a Large Language Model (LLM) with diverse tools and solvers through a conceptual interaction stack. The framework enables operations initiated by users, knowledge bases, or meta-learners to be processed by the LLM, which interfaces with specialized engines such as WolframAlpha and Whisper via conditions and constraints, enhancing the AI’s problem-solving capabilities.

```
1 :[Output]:
2 True
```

## E.2 DYNAMIC CASTING

By enabling sentence subtraction and dynamic casting within SymbolicAI, we utilize the generalization capability of NeSy engines to manipulate and refine text-based data, creating more meaningful and contextually relevant outcomes. The integration of dynamic casting with `Symbol` objects in our API allows the users to perform operations between `Symbol` objects and various data types, such as strings, integers, floats, lists, etc. without compromising on readability or simplicity.

```
1 s = symai.Symbol('Hello my enemy')
2 s - 'enemy' + 'friend'
```

```
1 :[Output]:
2 <class 'symai.expressions.Symbol'>(value=Hello my friend)
```

## E.3 SYMBOLS AND EMBEDDINGS

It is worth noting that encoding a complex object into a string sometimes precludes the object reconstitution. However, this concern does not substantially impede our methodology as we can employ approximations or proxy representations stored by the vector-valued property to effectively re-map objects. These representations are obtained through respective embedding models.

## E.4 TRANSLATION

In today’s increasingly interconnected world, translation between languages is fundamental, making it an essential feature.

```
1 s = symai.Symbol("Welcome to our tutorial.")
2 s.translate('German')
```

```
1 :[Output]:
2 <class 'symai.expressions.Symbol'>(value=Willkommen zu unserem Tutorial.)
```

## E.5 FILTERING, RANKING, EXTRACTION

Incorporating data-agnostic operations like filtering, ranking, and pattern extraction into our API allow the users to easily manipulate and analyze diverse data sets.

```
1 s = symai.Symbol(numpy.array([1, 2, 3, 4, 5, 6, 7]))
2 s.rank(measure='numerical', order='descending')

1 :[Output]:
2 <class 'symai.expressions.Symbol'>(value=['7', '6', '5', '4', '3', '2', '1'])
```

## E.6 IMPLICATIONS

One of the main objectives behind developing SymbolicAI was to facilitate reasoning capabilities in conjunction with the statistical inference inherent in LLMs. Consequently, we can carry out deductive reasoning operations through the Symbol objects. For instance, it is feasible to establish a series of operations with rules delineating the causal relationship between two symbols. The subsequent example illustrates the utilization of the & operator to compute the logical implication derived from the interaction of two distinct symbols.

```
1 s1 = symai.Symbol('The horn only sounds on Sundays.')
2 s2 = symai.Symbol('I hear the horn.')
3 s1 & s2

1 :[Output]:
2 <class 'symai.expressions.Symbol'>(value=It is Sunday.)
```

In the above example, the & operator overloads the and logical operator and extends its functionality. Furthermore, we can establish more sophisticated logical operators for and, or, and xor that can be grounded in formal proofs and utilize the NeSy engine to parse data structures before evaluating the expressions. This enables the definition of bespoke operations for executing intricate and robust logical operations, incorporating constraints to validate outcomes and guide the computation towards the desired behavior.

## E.7 CUSTOM OPERATIONS

The following example demonstrates how to define a custom == operation by overriding the `__eq__` method and providing a custom prompt object with a list of examples:

```
1 import symai
2
3 class Demo(symai.Symbol):
4     def __eq__(self, other) -> bool:
5         # define nested function
6         @symai.core.equals(examples=symai.Prompt([
7             "1 == 'ONE' =>True",
8             "'six' == 7 =>False",
9             "'Acht' == 'eight' =>True",
10            ...
11        ]))
12     def _func(_, other) -> bool: # [optional] cast return type (1. below)
13         return False # [optional] default behavior on failure (2. below)
14     return _func(self, other)
```

As illustrated in the example, this is also the method we used to implement basic operations in Symbol, namely by defining local functions that are then decorated with the respective operation decorator from the `symai.core.py` file. The `symai.core.py` is a collection of pre-defined operation decorators that can be quickly applied to any function. We use locally defined functions instead of directly decorating the main methods for two reasons:

1. We want to cast return types of the operation outcome to symbols or other derived classes thereof.
2. We do not necessarily want all of our operations to be sent to the NeSy engine and might need to implement a default behavior.

This is achieved through the `_sym_return_type` method, which can provide contextualized behavior based on the defined return type. More details can be found in the actual `Symbol` class.

In the context of LLMs, zero- and few-shot learning domains have emerged as essential techniques (Yao et al., 2023b; Shinn et al., 2023; Kim et al., 2023; Wei et al., 2022b; Lyu et al., 2023; Pitis et al., 2023; Madaan et al., 2023; Wang et al., 2022; Ye et al., 2023)<sup>7</sup> to enable models to generalize from limited training data and adapt to new tasks without requiring extensive retraining. This capability to learn and perform tasks with minimal examples is highly desirable, as it reduces the need for large labeled data sets and allows for faster deployment in new applications. In this section, we demonstrate how our Symbolic API incorporates Python decorators to define custom operations in the zero- and few-shot domains.

Consider the following example, where we define a custom operation to generate a random integer between 0 and 10 using the Symbolic API and Python decorators:

```

1 import symai
2
3 class Demo(symai.Symbol):
4     def __init__(self, value = '') -> None:
5         super().__init__(value)
6
7     @symai.core.zero_shot(prompt="Generate a random integer between 0 and 10.",
8                           constraints=[
9                               lambda x: x >= 0,
10                              lambda x: x <= 10
11                              ])
12     def get_random_int(self) -> int:
13         pass

```

In this example, the `@symai.core.zero_shot` decorator is used to define a custom operation that does not require any examples, as the prompt is expressive enough. The `zero_shot` decorator takes in two arguments: `prompt` and `constraints`. The prompt defines the conditioning for our desired operation behavior, while the constraints are used to validate the computed outcome, ensuring it meets our expectations. If the constraints are not fulfilled, the implementation would resort to the specified default implementation or the default value. If neither is provided, the Symbolic API raises a `ConstraintViolationException`. The return type in the example is defined as `int`. The resulting value from the wrapped function must be of type `int` because of the specific implementation of the auto-casting realized through `->`. If the cast fails, the Symbolic API raises a `ValueError`. If no return type is specified, the return type defaults to `Any`.

The `@symai.core.few_shot` decorator is a generalized version of `@symai.core.zero_shot` and is used to define custom operations requiring examples. The function signature of the `few_shot` decorator is as follows:

```

1 def few_shot(prompt: str,
2              examples: Prompt,
3              constraints: List[Callable] = [],
4              default: Any = None,
5              limit: int = 1,
6              pre_processor: Optional[List[PreProcessor]] = None,
7              post_processor: Optional[List[PostProcessor]] = None,
8              **wrp_kwargs):

```

The behavior of the `prompt` and `constraints` attributes is similar to the `zero_shot` decorator. The `examples` and `limit` arguments are new, with `examples` defining a list of instructions conditioning the NeSy engine, and `limit` specifying the maximum number of examples returned. The `pre_processor` and `post_processor` arguments accept lists of `PreProcessor` and `PostProcessor` objects, respectively, which are utilized to process the input before being fed into the NeSy engine and the output before being returned to the user. The `wrp_kwargs` argument passes additional arguments to the wrapped method, streamlining them towards the NeSy engine, or other engines.

## E.8 PROMPTING

In this section, we discuss the design of prompts and their role in shaping the behavior of operations. Prompts serve as containers for information necessary to define specific operations, and the `Prompt` class serves as the base class for

<sup>7</sup> This is by no means an exhaustive list, we only point the reader to some very influential and recent research.

all the other `Prompt` classes. Consider the following example, where we define a `Prompt` for comparing two values through the NeSy engine. In it, when the `<=` operation on two `Symbol` objects will be resolved, the NeSy engine evaluates them in the context of the `CompareValues` prompt.

```

1 class CompareValues(symai.Prompt):
2     def __init__(self) -> symai.Prompt:
3         super().__init__([
4             "4 > 88 =>False",
5             "-inf < 0 =>True",
6             "inf > 0 =>True",
7             "4 > 3 =>True",
8             "1 < 'four' =>True",
9             ...
10        ])

```

Evaluating a fuzzy comparison statement:

```

1 res = symai.Symbol(1) <= symai.Symbol('one')

```

Output of the evaluation:

```

1 :[Output]:
2 True

```

This evaluation returns `True`, as the fuzzy comparison operation conditions the engine to compare the two `Symbol` objects based on their semantic meaning. More generally, the semantics of `Symbol` operations may vary depending on the context hierarchy of the `Expression` class and the operations used. We used three main prompt designs: *Context-based Prompts*, *Operational Prompts*, and *Templates*. Prompts can be curated either through inheritance or composition. For instance, the *static context* can be defined by inheriting from the `Expression` class and overwriting the `static_context` property, while an `Operation` and `Template` prompt can be created by providing a `PreProcessor` to modify the input data.

We will now provide a more detailed explanation for each prompt design:

1. Context-based Prompts are considered optional and can be defined in a static manner, either by sub-classing the `Expression` class and overwriting the `static_context` property, or at runtime by updating the `dynamic_context` property or passing a payload kwargs to a method. Below is an example of payload kwargs through the method signature:

```

1 # creating a query to ask if an issue was resolve or not
2 s = symai.Symbol("<some-community-conversation>")
3 q = s.query("Was the issue resolved?")
4 # write manual condition to check if the issue was resolved
5 if 'not resolved' in q:
6     # do a new query but payload the previous query answer to the new query
7     s.query("What was the resolution?", payload=q)
8     ...
9 else:
10    pass # all good
11

```

Regardless of how the context is set, the contextualized prompt defines the desired behavior of `Expression` operations. For instance, if we want to operate in the context of a DSL without having to overwrite each base class method, we can utilize this approach<sup>8</sup>.

2. Operational Prompts define the behavior of an atomic operation and are therefore mandatory to express the nature of such an operation. For example, the `+` operation is used to add two `Symbol` objects together, and the `+` operation prompt explains its behavior. The `examples` kwargs provide another optional structure that conditions the NeSy engine with a set of instructions.
3. Template Prompts are optional and encapsulate the resulting prediction to enforce a specific format. For example, to generate HTML tags, we can utilize a curated `<html>...</html>` template. This template enforces the NeSy engine to begin the generation process already in the context of an HTML tags format and not produce irrelevant descriptions about its task.

<sup>8</sup> See more details in this [notebook](#).

## E.9 COMPLEX EXPRESSIONS

We will now attempt to obtain logical answers based on questions of the kind:

- A line parallel to  $y = 4x + 6$  passes through  $(5, 10)$ . What is the  $y$ -coordinate of the intercept?
- Bob has two sons, John and Jay. Jay has one brother and father. The father has two sons. Jay's brother has a brother and a father. Who is Jay's brother?
- Is 1000 bigger than 1063.472?

To solve these tasks, we would initially employ a series of operations to identify the most suitable engine for handling the specific requirements. Subsequently, we would prepare the input tailored to the selected engine.

```

1 val = "<one of the examples above>"
2
3 # First define a class that inherits from the \texttt{Expression} class
4 class ComplexExpression(symai.Expression):
5     # write a method that returns the causal evaluation
6     def causal_expression(self):
7         pass
8
9 # instantiate an object of the class
10 expr = ComplexExpression(val)
11 # set WolframAlpha as the main expression engine to use
12 expr.command(engines=['symbolic'], expression_engine='wolframalpha')
13 # evaluate the expression
14 res = expr.causal_expression()

```

The implementation of `causal_expression` could in principle look like this:

```

1 def causal_expression(self):
2     if self.isinstanceof('mathematics'):
3         # get the mathematical formula
4         formula = self.extract('mathematical formula')
5         # verify which problem type we have
6         if formula.isinstanceof('linear function'):
7             # prepare for wolframalpha
8             question = self.extract('question sentence')
9             req = question.extract('what is requested?')
10            # get coordinate point / could also ask for other points
11            x = self.extract('coordinate point (.,.)')
12            # concatenate to the question and formula
13            query = formula | f', point x = {x}' | f', solve {req}'
14            res = self.expression(query) # send prepared query to wolframalpha
15
16        elif formula.isinstanceof('number comparison'):
17            res = formula.expression() # send directly to wolframalpha
18
19        ... # more cases
20
21    elif self.isinstanceof('graph construction'):
22        sentences = self / '.' # first split into sentences
23        graph = {} # define graph
24        for s in sentences:
25            sym = symai.Symbol(s)
26            relations = sym.extract(
27                # extract and split by pipe
28                'connected entities (e.g. A has three B => A | A: three B)') / '|'
29        for r in relations:
30            ... # add relations and populate graph; reading suggestion
31
32        ... # more cases
33
34    return res

```

The aforementioned example demonstrates the utilization of the `causal_expression` method, which allows us to extract information that can be resolved either manually or through external solvers, say WolframAlpha. Initially, when utilizing the GPT-3 backend, we anticipated a significant engineering effort to develop such a complex expression, as the GPT-3 backend frequently struggled with accurate information extraction and comparison resolution. However, we remained confident in the field’s progress, specifically with fine-tuned models like RLHF ChatGPT. We were delighted to witness these challenges being further tackled through the latest GPT-4 model (OpenAI, 2023).

Furthermore, it is worth highlighting that, given sufficient data, we could refine methods for information extraction or knowledge graph construction from natural language, enabling more intricate reasoning tasks, such as those previously mentioned. We also direct readers to recent publications on Text-to-Graph translations, especially the very influential CycleGT (Guo et al., 2020). This approach allows us to answer queries by simply traversing the graph and extracting the required information.

Lastly, recent research (Kiciman et al., 2023; Ellis, 2023) has demonstrated that algorithms based on GPT-3.5 and GPT-4 establish new state-of-the-art accuracy on multiple causal benchmarks, while also exhibiting unique capabilities previously considered exclusive to humans, such as generating causal graphs and identifying background causal context from natural language. This points to the potential for LLMs to be used alongside existing causal methods as proxies for human domain knowledge, reducing human effort in setting up causal analyses and ultimately accelerating the widespread adoption of causal methods. Moreover, recent advances in LLMs have opened new frontiers for research, practice, and adoption of causal reasoning, transforming the way causal analysis is conducted and broadening the scope of applications for our framework.

One of the most prominent illustrations of this concept is exhibited by Word2Vec (Mikolov et al., 2013a). Word2Vec generates dense representations of words by training a shallow neural network to predict a word based on its neighboring words within a text corpus. These resulting vectors are extensively utilized in various natural language processing applications, including sentiment analysis, text classification, and clustering.

Drawing parallels with Word2Vec, our objective is to execute *contextualized* operations on different symbols. However, the key distinction lies in the fact that we operate within the natural language domain, as opposed to a vector space. Consequently, this grants us the capability to conduct arithmetic on words, sentences, paragraphs, and the like, while simultaneously validating the outcomes in a human-readable format.

The following example, we illustrate the methodology for evaluating such an expression through a string representation:

```
1 s = symai.Symbol('King - Man + Woman')
2 s.expression()

1 :[Output]:
2 <class 'symai.expressions.Symbol'> (value=Queen)
```

In contrast to the `Symbol` object, the `Expression` represents a non-terminal symbol. It allows for further evaluation and extends the `Symbol` class by overwriting the `__call__` method. It serves as the foundation for all other expressions and possesses additional capabilities, namely to `fetch` data from URLs, `search` the internet, or `open` files. These operations are intentionally separated from `Symbol`, as they do not utilize the `value` attribute of the `Symbol` class.

## E.10 COMPOSITION

### E.11 SEQUENCES

Sequences offer a multitude of advantages in the realm of `Expression` objects, as they facilitate the creation of more sophisticated structural configurations. By embodying the `Sequence` expression, multiple expressions can be effectively evaluated at runtime, thus enhancing the flexibility, modularity, and adaptability of the framework.

```
1 # first import all expressions
2 from symai.components import *
3 # define a sequence of expressions
4 Sequence(
5     Clean(),
6     Translate(),
7     Outline(),
8     Compose(),
9 )
```

## E.12 STREAMS

As demonstrated earlier, creating contextualized prompts refines the behavior of operations in the NeSy engine. However, this also consumes a considerable portion of the available context size. Given a limited context size, this constraint may pose challenges. Fortunately, the `Stream` processing expression offers a solution by opening a data stream and performing chunk-based operations on the input stream. `Stream` expressions can encapsulate other expressions. For instance, chunks can be managed through a `Sequence` expression, which permits multiple compositional operations sequentially. The example below illustrates the definition of a `Stream` expression:

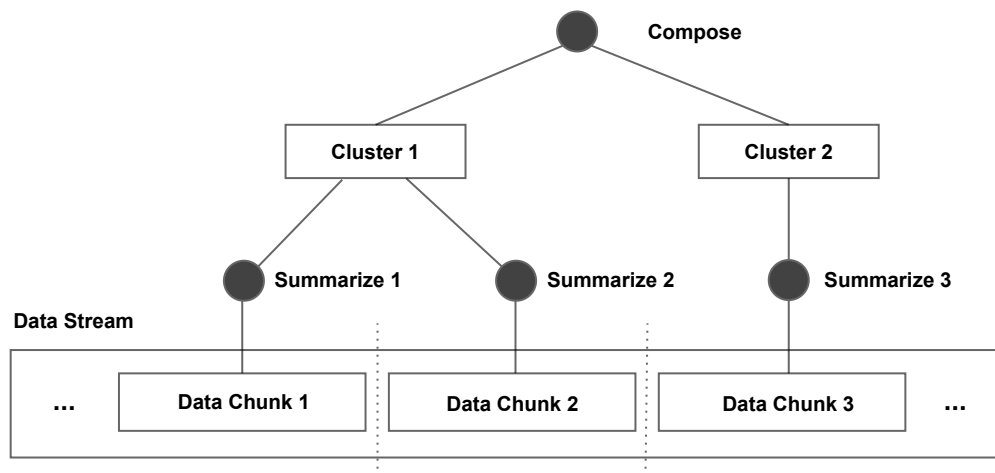
```

1 Stream(
2   Sequence(
3     Clean(),
4     Translate(),
5     Outline(),
6     Embed()
7   )
8 )

```

In this case, a stream is opened and a `Sequence` expression is passed, which cleans, translates, outlines, and embeds the input. Internally, the stream operation estimates the available model context size and segments the lengthy input text into smaller chunks transmitted to the inner expression. The returned object type is a generator.

The limitation of this approach is that the resulting chunks are processed independently, lacking shared context or information among them. To address this, the `Cluster` expression can be employed, merging the independent chunks based on their similarity, as it illustrated in Figure 12.



**Figure 12:** Stream processing expression in NeSy engine, illustrating data stream segmentation into chunks, each undergoing operations like cleaning, outlining, and embedding. The `Cluster` expression then merges chunks based on similarity, allowing contextually related information to be consolidated meaningfully. Node summaries are generated by extracting key labels from each cluster’s content, overcoming context size limitations and maintaining shared information among processed chunks.

By merging individual chunks by clustering their contents, contextually related information can be consolidated in a meaningful manner. Additionally, the clustered information can be labeled by streaming through each cluster’s content and extracting the most pertinent labels, yielding interpretable node summaries.

The complete example is depicted as follows:

```

1 stream = Stream(
2   Sequence(
3     Clean(),
4     Translate(),
5     Outline(),
6   )
7 )

```



```

8
9 s = symai.Symbol('<some long text>')
10 res = symai.Symbol(list(stream(s)))
11 expr = Cluster()
12 expr(res)

```

Subsequently, this process can be recursively repeated on each summary node to construct a hierarchical clustering structure. As each node represents a summarized subset of the original information, the summary can function as an index. The resulting tree can be utilized to navigate and retrieve the original information, transforming the large data stream problem into a search problem. Alternatively, vector-based similarity searches can be employed to identify similar nodes. For searching within a vector space, dedicated libraries such as Annoy (Spotify, 2017), Faiss (Johnson et al., 2019), or Milvus (Wang et al., 2021a) can be used.

In summary, `Stream` expressions offer the advantage of processing large data streams in a more efficient and organized manner, while also enabling the integration with other expressions like `Sequence` and `Cluster` expressions. These combinations allow for a more effective approach to handling context limitations, promoting the extraction of meaningful information and improving the overall performance of the framework.

### E.13 ERROR HANDLING, DEBUGGING, AND EXPLAINABILITY

Effective error handling and debugging are essential for ensuring the robustness and reliability of any software system, while explainability is essential for understanding the underlying behavior of the system, especially in the context of AI-driven applications. By developing a system that is both transparent and interpretable, we can gain valuable insights into the performance of the NeSy engines and identify potential areas for improvement.

### E.14 ERROR HANDLING

One of the fundamental aspects of the SymbolicAI API is being able to generate code. Consequently, errors may arise, and handling them contextually becomes vital. In pursuit of a self-evolving API, we introduce the `Try` expression, which includes built-in fallback statements and automatically retries execution after performing dedicated error analysis and correction. This expression analyzes both the input and the error, conditioning itself to resolve the error by modifying the original code<sup>9</sup>. If the fallback expression succeeds, the result is returned; otherwise, the process is repeated for the specified number of retries. If the maximum number of retries is reached without resolving the issue, the error is raised again.

Consider the example of executing previously generated code that contains a syntax error. By employing the `Execute` expression, we can evaluate the generated code, which takes a symbol and proceeds with the execution. Despite the initial failure, the `Try` expression resolves the syntactic error, returning the corrected and evaluated code:

```

1 expr = Try(expr=Execute())
2 s = symai.Symbol('a = int("3,")') # some code with a syntax error
3 expr(s)

1 :Output:
2 a = 3

```

While not all errors can be resolved as easily as the demonstrated syntax error, we continue to explore more sophisticated error handling mechanisms, including streams and clustering to address errors in a hierarchical and contextual manner.

## F EVALUATION DETAILS

In a computational graph, the VERTEX score compares the distribution of the generated model answer at each node against a reference distribution by sampling multiple valid trajectories at each node for the reference distribution and accounting for randomness through some predefined random trajectories. For instance, one of the predefined random trajectories in our benchmark was the string of ASCII characters which are considered printable, namely 0123456789abcdefghijklmnopqrstuvwxyzABCDEFGHIJKLMNOPQRSTUVWXYZ! "#\$%&'()\*+,-./:;<=>?@[ ]^\_`{|}~. Moreover, the VERTEX score is particularly well suited for the evaluation of multi-step

<sup>9</sup> This is similar to the recently released [Auto-GPT](#) application.

workflows and in contexts where the solution space is or is expected to be diverse. We will now proceed by describing in detail the tasks that we defined in our benchmark.

### F.1 ASSOCIATIVE PREDICTION

We defined a total of 15 tasks involving in-context associations between two `Symbol` instances. SymbolicAI’s overloaded operators rely on predefined pseudo-grammars, as described in Section 4, that augment the operators with few-shot examples. For instance, the overloaded operator `+` utilized between two `Symbol` instances provides few-shot examples how to resolve additions with various data types:

```

1 "'1' + 2 =>3",
2 "17 + 'pi' =>20.1415926535...",
3 "7.2 + 'five' =>12.2",
4 "True + 0 => False",
5 "False + 'True' =>False",
6 "['a', 'b'] + ['c', 'd'] =>['a', 'b', 'c', 'd']",
7 "['apple'] + 'banana' =>['apple', 'banana']",
8 "'Zero' + 1 =>1",
9 "'One' + 'Two' =>3",
10 "'Three' + 4 =>7",
11 "'a + b' + 'c + d' =>a + b + c + d",
12 ...

```

Consequently, we can now test if the models can solve the addition between `Symbol("two hundred and thirty four")` and `Symbol(7000)`.

### F.2 MULTI-MODAL BINDING

We perform transformations between multiple modalities through language-based representations. Therefore, we need to evaluate the model’s proficiency in tool utilization, classification and routing of requests to relevant modules. We define a multi-modal `Expression` to detect the category of a task based on its content and to forward the task to the appropriate tool. The expression creates interfaces to tools like WolframAlpha for mathematical expressions, Selenium for website content scraping, SerpApi for search queries, and APILayer for optical character recognition. Each of the five tests aims to evaluate the appropriate handling of a specific type of input by the multi-modal `Expression` type, such as processing a website URL for scraping, interpreting a search engine query, testing if two vectors are linearly independent, comparing large numbers, and extracting text from an image. The following example shows the `MultiModalExpression` implementation of the `forward` function that uses `isinstanceof` operator on its own context to determine its current expression value and select the sub-routine that can evaluate the request.

```

1 class MultiModalExpression(Expression):
2     def forward(self, ...):
3         formula = self.extract('mathematical formula')
4         ...
5         if self.isinstanceof(LINEAR_ALGEBRA):
6             ...
7             res = self.solver(formula)
8             res = res.query('write a one sentence summary of the answer')
9             ...
10        elif self.isinstanceof(NUMBER_COMPARISON):
11            res = self.solver(formula) # send directly to wolframalpha
12        else:
13            ...
14
15 query = Symbol("is 100044347 bigger than 129981063.472?")
16 expr = MultiModalExpression(query)
17 res = expr(...)

```

### F.3 PROGRAM SYNTHESIS

We designed three separate tests related to program synthesis, where each task assesses the ability of the models to generate and execute code based on natural language instructions or provided templates:

- 1) The first task involves reading a LaTeX table template and data, then generating a function to populate the table with the given data.
- 2) The second task tests the automatic code generation for API calls by fetching data from a specified URL and extracting specific information from the retrieved content.
- 3) The third task evaluates the ability to construct a custom Expression that processes a Symbol through a specific Function component from the SymbolicAI package.

Each of the three tests follows a similar pattern, where the generated code is scored based on its similarity to valid references and normalized with random samples. Figure 13 shows possible samples from the third task category.



**Figure 13:** Samples for constructing custom expressions. The expected references are human-generated and are distributed according to the solution space. The set of valid references can be as well human-generated. The set of invalid (or undesirable) samples represent the subset we are the least interested in. Lastly, note that in general, the references need not be human-generated; any combination of human- and machine-generated references is possible. For instance, if synthetic data is utilized to distill knowledge from a larger model to a smaller model, the valid references and the expected references can be sampled from the larger model.

#### F.4 LOGICAL COMPONENTS

We designed six tests to assess the logical capabilities of the candidate models and group them as follows.

- 1) We utilize the Python library SymPy for symbolic mathematics to create the mathematical expression  $ax + bx - cx - ay - by + cy + d$ . The task for the model is then to factorize the expression and extract all unique symbols as a list.

```
1 import sympy as sym
2
3 ...
4 a, b, c, d, x, y = sym.symbols('a, b, c, d, x, y')
```

```

5 expr = a * x + b * x - c * x - a * y - b * y + c * y + d
6 # validate with sympy
7 fact = sym.collect(expr, d, func=sym.factor)
8 # model based factorization
9 func = Factorization('Factorize d from the expression such that your final start with:
    `d + (...):')
10 res = func(expr)
11 # compare res with fact
12 ...

```

2) Three tasks evaluate a models' capability to resolve the logical operations AND, OR, and XOR. For instance, the test for logical AND combines the symbols `Symbol("The horn only sounds on Sundays")` and `Symbol("I hear the horn")` and compares the answer against the human-generated references *"The horn only sounds on Sundays and I hear the horn."* and *"Since I hear the horn it is Sunday."* Since there is a large number of possible solutions, there is high variability in the solution space. Each model might prefer a different solution.

The following snippet shows how one can define a custom primitive class (`CustomLogicPrimitive`) for logic operators. The `__or__` function gets overloaded and uses the built-in `logic` decorator from the `core` package to create a local function that evaluates two `Symbol` instances.

```

1 from symai import core
2 from symai.ops.primitives import Primitive
3
4 class CustomLogicPrimitive(Primitive):
5     def __or__(self, other: Any) -> Any:
6         @core.logic(operator='or')
7         def _func(_, a: str, b: str):
8             pass # could impl. a fallback behavior here
9             return self._to_symbol(_func(self, other))
10    ...
11
12 subject = 'cat'
13 res = (Symbol(f'The {subject} has whiskers.', primitives=CustomLogicPrimitive) | \
14        Symbol(f'The {subject} has a tail.', primitives=CustomLogicPrimitive))

```

3) For another task we use a custom `Expression` that defines a DSL syntax and semantic structure. We use this `Expression` to extract higher-order logic expressions from a natural language statement, namely the puzzle 'Who is Jay's brother?'<sup>10</sup>, that preserves the original relationships. The following is a DSL snippet of the 'Who is Jay's brother?' puzzle:

```

1 // Query
2 IsBrotherOf(jay, john, bob) <- BrotherOf(jay, john) AND FatherOf(bob, jay) AND
    FatherOf(bob, john);
3
4 // Facts
5 BrotherOf(x, y) <- HAS(x, brother) AND HAS(y, brother) AND Sibling(x, y);
6 FatherOf(x, y) <- HAS(x, son) AND ParentOf(x, y);
7 ParentOf(x, y) <- IS(x, parent) AND IS(y, child);
8 Sibling(x, y) <- IS(x, father) AND IS(y, father) OR IS(x, mother) AND IS(y, mother);
9
10 ...

```

4) For the final task, we again use the puzzle 'Who is Jay's brother?' to evaluate a models' capability for complex conversions. We use the Z3 theorem prover (Moura & Bjørner, 2008) to solve the 'Who is Jay's brother' puzzle conditioned on the Z3 solvers' solution to Einsteins' famous puzzle 'Who owns the fish?'. The task involves an indirect translation from natural language to executable code by the Z3 solver; the solution to Einstein's puzzle acts as a form of self-contained "documentation" for how the Z3 solver should be utilized. The test constructs a template, which includes the task instructions, puzzle statement, and reference to the Einstein's puzzle solution. The models are then asked to analyze the given problem and solution format and create a Python function with Z3 syntax that can solve the 'Who is Jay's brother?' puzzle. The dynamically generated code is executed within the test environment utilizing Python's `exec` function. We check the access to the Z3 solver and run the generated `solve_puzzle`

<sup>10</sup> Bob has two sons, John and Jay. Jay has one brother and father. The father has two sons. Jay's brother has a brother and a father. Who is Jay's brother?

function supposed to contain the logic to solve the puzzle. Once executed, the assembled Z3 logical clauses are processed by the solver, which verifies that the set of constraints is satisfiable. If so, the model generated by the solver is queried for the puzzle’s solution and scored using our VERTEX score. The following is an example output from the Z3 representation of the solution to ‘Who is Jay’s brother?’ puzzle:

```

1 from z3 import Solver, Bool, And, Not, Const, BoolSort, EnumSort, Function, IntSort
2
3 def solve_puzzle(S: Solver) -> Const:
4     # Define the enumeration sort for the individuals
5     Person, (BobE, JohnE, JayE, JaysBrotherE, FatherE) = EnumSort('Person', ['Bob', '
        John', 'Jay', 'JaysBrother', 'Father'])
6
7     # Define a function from boolean to persons (for brother status)
8     is_brother = Function('is_brother', Person, BoolSort())
9
10    # Define the relationships
11    S.add(is_brother(JohnE) == True) # John is a brother
12    S.add(is_brother(JayE) == True) # Jay is a brother
13
14    ...
15
16    return query

```

## F.5 HIERARCHICAL COMPUTATIONAL GRAPHS

In this section we extend on the hierarchical computational graphs section.

**Research Paper Draft** The following example defines a `Paper` expression that takes in a sequence of expressions which are sequentially executed. The `Method` expression contains a `Source` expression which points to the actual human-written method. The `Method` expression acts as the root node that bootstraps the generation process. The `RelatedWork` expression contains a sequence of `Cite` expressions which are executed in parallel and utilized to define the context of the related work section. The `Abstract` and `Title` expressions get executed last because they require all the previous information to be available in their respective contexts. Each expression in the sequence of expressions from `Paper` takes in the context of its predecessors. All expressions also link their results to a global linker object which is utilized after the execution to retrieve individual results from the nodes of the expression’s computational graph. Each node was evaluated against its corresponding reference, all references representing actual sections from this research paper. The samples for each node were generated with a separate model (Claude 2) that was not part of this evaluation. In Figure 5 we show the resulting computational graph of the `Paper` expression.

```

1 # define the computational graph
2 expression = Paper(
3     Method(
4         # link to original code base where the main method is defined
5         Source(file_link='/path/to/.../file'),
6     ),
7     # gather resources and write the related work
8     RelatedWork(
9         Cite(file_link='/path/to/.../file'),
10        Cite(file_link='/path/to/.../file'),
11        ...
12    ),
13    # write the abstract and title
14    Abstract(),
15    Title(),
16 )
17 # run the graph
18 paper_result = expression('Write a scientific paper')
19 # access linker to retrieve the results from the method expression
20 method = expr.linker.find('Method')
21 # print result of the method expression
22 print(method)

```

## F.6 CAVEATS

One may have noticed that there are cases in which there is no need to sample multiple samples as there is only one expected answer, for instance, if we need to extract a specific number from a string and cast it to an integer. If the extraction process appends additional characters other than the number, the casting will fail. In such cases, the VERTEX score simply defaults to the chosen similarity measure, registering and penalizing any deviation from the expected answer.